

10/033-167

(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
11 January 2001 (11.01.2001)

PCT

(10) International Publication Number
WO 01/02568 A2

(51) International Patent Classification⁷: C12N 15/12,
15/55, 15/54, 15/61, C07K 14/47, C12N 9/64, 9/12, 9/90,
C12Q 1/68, C12N 15/11, C07K 16/18, 16/40, G01N
33/566, A61K 38/00

(21) International Application Number: PCT/US00/18374

(22) International Filing Date: 30 June 2000 (30.06.2000)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:
60/142,310 2 July 1999 (02.07.1999) US
60/142,311 2 July 1999 (02.07.1999) US

(71) Applicants: CHIRON CORPORATION [US/US]; 4560
Horton Street, Emeryville, CA 94608 (US). HYSEQ,
INC. [US/US]; 675 Almanor Avenue, Sunnyvale, CA
94086 (US).

(72) Inventors: WILLIAMS, Lewis, T.; Chiron Corporation,
P.O. Box 8097, Emeryville, CA 94662-8097 (US). ES-
COBEDO, Jaime; Chiron Corporation, P.O. Box 8097,
Emeryville, CA 94662-8097 (US). INNIS, Michael, A.;
Chiron Corporation, P.O. Box 8097, Emeryville, CA
94662-8097 (US). GARCIA, Pablo, Dominguez; Chiron
Corporation, P.O. Box 8097, Emeryville, CA 94662-8097
(US). KLINGER, Julie; Chiron Corporation, P.O. Box
8097, Emeryville, CA 94662-8097 (US). KASSAM,
Altaf; Chiron Corporation, P.O. Box 8097, Emeryville,
CA 94662-8097 (US). REINHARD, Christoph; Chiron
Corporation, P.O. Box 8097, Emeryville, CA 94662-8097
(US). RANDAZZO, Filippo; Chiron Corporation,
P.O. Box 8097, Emeryville, CA 94662-8097 (US).
KENNEDY, Guilina, C.; Chiron Corporation, P.O. Box
8097, Emeryville, CA 94662-8097 (US). POT, David;
Chiron Corporation, P.O. Box 8097, Emeryville, CA
94662-8097 (US). LAMSON, George; Chiron Corpora-
tion, P.O. Box 8097, Emeryville, CA 94662-8097 (US).

DRMANAC, Radoje; 675 Almanor Avenue, Sunnyvale,
CA 94086 (US). CRKENJAKOV, Radomir; 675 Al-
manor Avenue, Sunnyvale, CA 94086 (US). DRMANAC,
Snezana; 675 Almanor Avenue, Sunnyvale, CA 94086
(US). DICKSON, Mark; 675 Almanor Avenue, Sun-
nyvale, CA 94086 (US). LABAT, Ivan; 675 Almanor
Avenue, Sunnyvale, CA 94086 (US). LESHKOWITZ,
Dena; 675 Almanor Avenue, Sunnyvale, CA 94086 (US).
KITA, David; 675 Almanor Avenue, Sunnyvale, CA
94086 (US). GARCIA, Veronica; 675 Almanor Avenue,
Sunnyvale, CA 94086 (US). JONES, Lee, William; 675
Almanor Avenue, Sunnyvale, CA 94086 (US). STRA-
CHE-CRAIN, Birgit; 675 Almanor Avenue, Sunnyvale,
CA 94086 (US).

(74) Agents: BLACKBURN, Robert, P.; Chiron Corporation,
4560 Horton Street, Emeryville, CA 94608-2916 et al.
(US).

(81) Designated States (*national*): AE, AG, AL, AM, AT, AU,
AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CR, CU, CZ,
DE, DK, DM, DZ, EE, ES, FI, GB, GD, GE, GH, GM, HR,
HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR,
LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ,
NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM,
TR, TT, TZ, UA, UG, UZ, VN, YU, ZA, ZW.

(84) Designated States (*regional*): ARIPO patent (GH, GM,
KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian
patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European
patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE,
IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG,
CI, CM, GA, GN, GW, ML, MR, NE, SN, TD, TG).

Published:

— Without international search report and to be republished
upon receipt of that report.

For two-letter codes and other abbreviations, refer to the "Guid-
ance Notes on Codes and Abbreviations" appearing at the begin-
ning of each regular issue of the PCT Gazette.

WO 01/02568 A2

(54) Title: NOVEL HUMAN GENES AND GENE EXPRESSION PRODUCTS

(57) Abstract: The invention provides novel polynucleotides. The invention further provides novel members of protein families, and polynucleotides that are differentially expressed in cancer cells relative to normal cells, and in metastatic cancer cells relative to normal cells or non-metastatic cancer cells.

NOVEL HUMAN GENES AND GENE EXPRESSION PRODUCTS

FIELD OF THE INVENTION

The present invention relates to novel polynucleotides of human origin and the encoded gene products.

5 BACKGROUND OF THE INVENTION

Identification of novel polynucleotides, particularly those that encode an expressed gene product, is important in the advancement of drug discovery, diagnostic technologies, and the understanding of the progression and nature of complex diseases such as cancer. Identification of genes expressed in different cell types isolated from
10 sources that differ in disease state or stage, developmental stage, exposure to various environmental factors, the tissue of origin, the species from which the tissue was isolated, and the like is key to identifying the genetic factors that are responsible for the phenotypes associated with these various differences.

This invention provides novel human polynucleotides, the polypeptides
15 encoded by these polynucleotides, and the genes and proteins corresponding to these novel polynucleotides.

SUMMARY OF THE INVENTION

This invention relates to novel human polynucleotides and variants thereof, their encoded polypeptides and variants thereof, to genes corresponding to these
20 polynucleotides and to proteins expressed by the genes. The invention also relates to diagnostics and therapeutics comprising such novel human polynucleotides, their corresponding genes or gene products, including probes, antisense nucleotides, and antibodies. The polynucleotides of the invention correspond to a polynucleotide comprising the sequence information of at least one of SEQ ID NOs: 1-3351.

25 Various aspects and embodiments of the invention will be readily apparent to the ordinarily skilled artisan upon reading the description provided herein.

DETAILED DESCRIPTION OF THE INVENTION

The invention relates to polynucleotides comprising the disclosed nucleotide sequences, to full length cDNA, mRNA genomic sequences, and genes

corresponding to these sequences and degenerate variants thereof, and to polypeptides encoded by the polynucleotides of the invention and polypeptide variants.

Polypeptide variants differ from wild type protein in having one or more amino acid substitutions that either enhance, add, or diminish a biological activity of the wild type protein.

Six of the polypeptides disclosed herein encode new members of the MKK kinase family; the coding region is found within the nucleotide region in parentheses: SEQ ID NO:29 (nucleotides 295-421); SEQ ID NO:31 (298-397); SEQ ID NO:196 (37-322); SEQ ID NO:3175 (nucleotides 14-164); SEQ ID NO:3190 (229-390); and SEQ ID NO:3281 (15-182). Twenty-four of the polypeptides encode new members of the family of transcription factor proteins having a basic region plus leucine zipper: SEQ ID NO:410 (42-191); SEQ ID NO:552 (116-288); SEQ ID NO:768 (116-288); SEQ ID NO:822 (108-262); SEQ ID NO:836 (158-353); SEQ ID NO:1288 (73-234); SEQ ID NO:1365 (69-257); SEQ ID NO:1540 (289-471); SEQ ID NO:1549 (200-391); SEQ ID NO:1556 (163-354); SEQ ID NO:1557 (207-398); SEQ ID NO:1563 (107-298); SEQ ID NO:1622 (180-365); SEQ ID NO:1630 (100-291); SEQ ID NO:1704 (184-372); SEQ ID NO:1808 (36-161); SEQ ID NO:1454 (49-209); SEQ ID NO:2363 (48-211); SEQ ID NO:2424 (43-194); SEQ ID NO:3147 (190-369); SEQ ID NO:3152 (129-320); SEQ ID NO:3158 (167-334); and SEQ ID NO:3208 (34-256).

SEQ ID NOs:186 (175-395); 2591 (60-165); 3307 (43-321); and 3339 (94-342) encode polypeptides having an SH2 domain, and SEQ ID NOs:234 (23-121), 1832 (18-173), and 1835 (57-206) encode polypeptides having an SH3 domain. Nine polypeptides encode new members of the family of proteins having Ank repeat regions: SEQ ID NO:187 (358-432); SEQ ID NO:1268 (238-315); SEQ ID NO:1804 (301-378); SEQ ID NO:1819 (278-355); SEQ ID NO:1839 (224-307); SEQ ID NO:1830 (184-267); SEQ ID NO:2562 (18-101); SEQ ID NO:3015 (131-214); and SEQ ID NO:3267 (97-180).

The following eleven polynucleotides encode polypeptides having a C2H2 type zinc finger: SEQ ID NOs:308 (110-172); 807 (339-392); 1324 (294-356); 1503 (154-216); 1527 (156-212); 1674 (196-258); 1779 (64-126); 1801 (295-351); 3081 (190-252); 3193 (293-355); and 3306 (161-223). Eight polynucleotides encode polypeptides of the family of ATPases: SEQ ID NOs:431 (71-428); 639 (157-561); 2135 (2-401); 2684 (9-461); 2859 (100-320); 3178 (45-386); 3197 (281-343) and 3266 (8-139). Polypeptides having a fibronectin type III domain are encoded by SEQ ID NO:746 (209-427) and 1192 (186-416). Polypeptides having an EF-hand domain are encoded by SEQ ID NO:820 (341-

406); 1755 (281-367) and 3285(16-102). Six polypeptides of the protein kinase family are encoded by SEQ ID NOs:1157 (41-444); 1478 (54-437), 1496 (241-520); 2286 (12-182); 2969 (5-387); and 3190 (118-390).

LIM domain-containing polypeptides are encoded by SEQ ID NO:1269 (79-240); 1309 (248-404); 1360 (222-377); and 1386 (243-398). Two polypeptides of the family having a C2 domain (protein kinase C-like) are encoded by SEQ ID NO:1325 (1-234) and 2282(183-353). Polypeptides having a WD domain, G-beta repeat motif are encoded by SEQ ID NOs:1336 (66-164); 1380 (42-140); 1711 (263-361); 1762 (236-334); 1909 (160-258); 2218 (127-225); 3047 (191-292); 3108 (275-367) and 3292 (208-300).

SEQ ID NO:1410 (222-350) encodes a member of the trypsin family. SEQ ID NOs:1417 (8-354); 2281 (20-387) and 2310 (20-371) encode members of the protein tyrosine phosphatase family. SEQ ID NOs:1464 (4-180) and 1514 (2-252) encode members of the family having an RNA recognition motif (also known as RRM, RBD, or RNP domain). SEQ ID NOs:1496 (241-520) and 3297(7-153) encode helicases having a conserved C-terminal domain. SEQ ID NO:1538 (9-635) encodes a member of the wnt family of developmental signaling proteins.

Three polynucleotides encode polypeptides having a homeobox domain: SEQ ID NOs:1676 (9-86); 1820 (123-299); and 1821 (127-303). A novel thioredoxin is encoded by SEQ ID NO:1677 (316-369). Two novel members of the ras family are encoded by SEQ ID NO:1688(109-410) and 3258(138-394). A novel polypeptide having a phosphatidylinositol-specific phospholipase C Y-domain is encoded by SEQ ID NO:1707 (92-439). A novel serine carboxypeptidase is encoded by SEQ ID NO:1744 (238-433). A novel polypeptide having N-terminal homology in the Ets domain is encoded by SEQ ID NO:1811 (184-315). A novel polypeptide having a bromodomain is encoded by SEQ ID NO:1814 (127-294). A novel polypeptide having a double-stranded RNA binding motif is encoded by SEQ ID NO:1818 (9-146). A novel polypeptide having a G-protein alpha subunit is encoded by SEQ ID NO:1846 (12-398).

SEQ ID NOs:1911 (35-151) and 1980 (60-197) encode polypeptides having a C3HC4 type zinc finger domain (RING finger). SEQ ID NO:2065 (253-306) encodes a polypeptide having a CCHC zinc finger domain. SEQ ID NO:2216 (90-179) encodes a polypeptide having a WW/rsp5/WWP domain. SEQ ID NO:2428 (25-350) encodes a polypeptide member of the dual specificity phosphatase family, having a catalytic domain.

SEQ ID NOs:2577 (0-311); 3183 (14-215); and 3195 (0-215) encode members of the 4 transmembrane segment integral membrane protein family. SEQ ID

NOs:2826 (116-400) and 2871 (198-392) encode polypeptides of the DEAD and DEAH box helicase family. SEQ ID NO:2944 (18-281) encodes a polypeptide having a calpain large subunit, domain III.

- 5 SEQ ID NO:3274 (11-187) encodes a eukaryotic transcription factor with a fork head domain. SEQ ID NO:3345 (65-271) encodes a polypeptide having a PDZ domain, and SEQ ID NO:3351 (124-270) encodes a polypeptide in the family of phorbol esters/glycerol binding proteins.

- Described below are polynucleotide compositions encompassed by the invention, methods for obtaining cDNA or genomic DNA encoding a full-length gene product, expression of these polynucleotides and genes, identification of structural motifs of the polynucleotides and genes, identification of the function of a gene product encoded by a gene corresponding to a polynucleotide of the invention, use of the provided polynucleotides as probes and in mapping and in tissue profiling, use of the corresponding polypeptides and other gene products to raise antibodies, and use of the polynucleotides and their encoded gene products for therapeutic and diagnostic purposes.
- 10
15

Polynucleotide Compositions

- The scope of the invention with respect to polynucleotide compositions includes, but is not necessarily limited to, polynucleotides having a sequence set forth in any one of SEQ ID NOs:1-3351; polynucleotides obtained from the biological materials described herein or other biological sources (particularly human sources) by hybridization under stringent conditions (particularly conditions of high stringency); genes corresponding to the provided polynucleotides; variants of the provided polynucleotides and their corresponding genes, particularly those variants that retain a biological activity of the encoded gene product (e.g., a biological activity ascribed to a gene product corresponding to the provided polynucleotides as a result of the assignment of the gene product to a protein family(ies) and/or identification of a functional domain present in the gene product). Other nucleic acid compositions contemplated by and within the scope of the present invention will be readily apparent to one of ordinary skill in the art when provided with the disclosure here.
- 20
25
30
- "Polynucleotide" and "nucleic acid" as used herein with reference to nucleic acids of the composition is not intended to be limiting as to the length or structure of the nucleic acid unless specifically indicated.

The invention features polynucleotides that are expressed in human tissue, specifically human colon, breast, and/or lung tissue. Novel nucleic acid

compositions of the invention comprise a sequence set forth in any one of SEQ ID NOs:1-3351 or an identifying sequence thereof. An "identifying sequence" is a contiguous sequence of residues at least about 10 nt to about 20 nt in length, usually at least about 50 nt to about 100 nt in length, that uniquely identifies a polynucleotide sequence, *e.g.*, exhibits less than 90%, usually less than about 80% to about 85% sequence identity to any contiguous nucleotide sequence of more than about 20 nt. Thus, the subject novel nucleic acid compositions include full length cDNAs or mRNAs that encompass an identifying sequence of contiguous nucleotides from any one of SEQ ID NOs:1-3351.

10 The polynucleotides of the invention also include polynucleotides having sequence similarity or sequence identity. Nucleic acids having sequence similarity are detected by hybridization under low stringency conditions, for example, at 50°C and 10XSSC (0.9 M saline/0.09 M sodium citrate) and remain bound when subjected to washing at 55°C in 1XSSC. Sequence identity can be determined by hybridization
15 under stringent conditions, for example, at 50°C or higher and 0.1XSSC (9 mM saline/0.9 mM sodium citrate). Hybridization methods and conditions are well known in the art, see, *e.g.*, U.S. Patent No. 5,707,829. Nucleic acids that are substantially identical to the provided polynucleotide sequences, *e.g.*, allelic variants, genetically altered versions of the gene, *etc.*, bind to the provided polynucleotide sequences (SEQ
20 ID NOs:1-3351) under stringent hybridization conditions. By using probes, particularly labeled probes of DNA sequences, one can isolate homologous or related genes. The source of homologous genes can be any species, *e.g.*, primate species, particularly human; rodents, such as rats and mice; canines, felines, bovines, ovines, equines, yeast, nematodes, *etc.*

25 Preferably, hybridization is performed using at least 15 contiguous nucleotides (nt) of at least one of SEQ ID NOs:1-3351. That is, when at least 15 contiguous nt of one of the disclosed SEQ ID NOs. is used as a probe, the probe will preferentially hybridize with a nucleic acid comprising the complementary sequence, allowing the identification and retrieval of the nucleic acids that uniquely hybridize to
30 the selected probe. Probes from more than one SEQ ID NO. can hybridize with the same nucleic acid if the cDNA from which they were derived corresponds to one mRNA. Probes of more than 15 nt can be used, *e.g.*, probes of from about 18 nt to about 100 nt, but 15 nt represents sufficient sequence for unique identification.

The polynucleotides of the invention also include naturally occurring
35 variants of the nucleotide sequences (*e.g.*, degenerate variants, allelic variants).

5 Variants of the polynucleotides of the invention are identified by hybridization of putative variants with nucleotide sequences disclosed herein, preferably by hybridization under stringent conditions. For example, by using appropriate wash conditions, variants of the polynucleotides of the invention can be identified where the allelic variant exhibits at most about 25-30% base pair (bp) mismatches relative to the selected polynucleotide probe. In general, allelic variants contain 15-25% bp mismatches, and can contain as little as even 5-15%, or 2-5%, or 1-2% bp mismatches, as well as a single bp mismatch.

10 The invention also encompasses homologs corresponding to the polynucleotides of SEQ ID NOs:1-3351, where the source of homologous genes can be any mammalian species, *e.g.*, primate species, particularly human; rodents, such as rats; canines, felines, bovines, ovines, equines, yeast, nematodes, *etc.* Between mammalian species, *e.g.*, human and mouse, homologs generally have substantial sequence similarity, *e.g.*, at least 75% sequence identity, usually at least 90%, more usually at 15 least 95% between nucleotide sequences. Sequence similarity is calculated based on a reference sequence, which may be a subset of a larger sequence, such as a conserved motif, coding region, flanking region, *etc.* A reference sequence will usually be at least about 18 contiguous nt long, more usually at least about 30 nt long, and may extend to the complete sequence that is being compared. Algorithms for sequence analysis are 20 known in the art, such as BLAST, described in Altschul et al., *J. Mol. Biol.* (1990) 215:403-10.

In general, variants of the invention have a sequence identity greater than at least about 65%, preferably at least about 75%, more preferably at least about 85%, and can be greater than at least about 90%, 91%, 92%, 93%, 94%, 95%, or 96%, most 25 preferably 97%, 98% or 99%. For the purposes of this invention, a preferred method of calculating percent identity is the Smith-Waterman algorithm, using the following. Global DNA sequence identity must be greater than 65% as determined by the Smith-Waterman homology search algorithm as implemented in MPSRCH program (Oxford Molecular) using an affine gap search with the following search parameters: gap open 30 penalty, 12; and gap extension penalty, 1.

The subject nucleic acids can be cDNAs or genomic DNAs, as well as fragments thereof, particularly fragments that encode a biologically active gene product and/or are useful in the methods disclosed herein (*e.g.*, in diagnosis, as a unique identifier of a differentially expressed gene of interest, *etc.*). The term "cDNA" as used 35 herein is intended to include all nucleic acids that share the arrangement of sequence

elements found in native mature mRNA species, where sequence elements are exons and 3' and 5' non-coding regions. Normally mRNA species have contiguous exons, with the intervening introns, when present, being removed by nuclear RNA splicing, to create a continuous open reading frame encoding a polypeptide of the invention.

5 A genomic sequence of interest comprises the nucleic acid present between the initiation codon and the stop codon, as defined in the listed sequences, including all of the introns that are normally present in a native chromosome. It can further include the 3' and 5' untranslated regions found in the mature mRNA. It can further include specific transcriptional and translational regulatory sequences, such as
10 promoters, enhancers, *etc.*, including about 1 kb, but possibly more, of flanking genomic DNA at either the 5' and 3' end of the transcribed region. The genomic DNA can be isolated as a fragment of 100 kbp or smaller; and substantially free of flanking chromosomal sequence. The genomic DNA flanking the coding region, either 3' and 5', or internal regulatory sequences as sometimes found in introns, contains sequences
15 required for proper tissue, stage-specific, or disease-state specific expression.

The nucleic acid compositions of the subject invention can encode all or a part of the subject polypeptides. Double or single stranded fragments can be obtained from the DNA sequence by chemically synthesizing oligonucleotides in accordance with conventional methods, by restriction enzyme digestion, by PCR amplification, *etc.*
20 Isolated polynucleotides and polynucleotide fragments of the invention comprise at least about 10, about 15, about 20, about 35, about 50, about 100, about 150 to about 200, about 250 to about 300, or about 350 contiguous nt selected from the polynucleotide sequences as shown in SEQ ID NOs:1-3351. The fragments also include those of lengths intermediate to the specifically mentioned lengths, such as 35,
25 36, 37, 38, 39, *etc.*; 150, 151, 152, 153, 154, *etc.* For the most part, fragments will be of at least 15 nt, usually at least 18 nt or 25 nt, and up to at least about 50 contiguous nt in length or more. In a preferred embodiment, the polynucleotide molecules comprise a contiguous sequence of at least 12 nt selected from the group consisting of the polynucleotides shown in SEQ ID NOs:1-3351.

30 Probes specific to the polynucleotides of the invention can be generated using the polynucleotide sequences disclosed in SEQ ID NOs:1-3351. The probes are preferably at least about a 12, 15, 16, 18, 20, 22, 24, or 25 nt fragment of a corresponding contiguous sequence of SEQ ID NOs:1-3351, and can be less than 2, 1, 0.5, 0.1, or 0.05 kb in length. The probes can be synthesized chemically or can be
35 generated from longer polynucleotides using restriction enzymes. The probes can be

labeled, for example, with a radioactive, biotinylated, or fluorescent tag. Preferably, probes are designed based upon an identifying sequence of a polynucleotide of one of SEQ ID NOs:1-3351. More preferably, probes are designed based on a contiguous sequence of one of the subject polynucleotides that remain unmasked following
5 application of a masking program for masking low complexity (*e.g.*, XBLAST) to the sequence., *i.e.*, one would select an unmasked region, as indicated by the polynucleotides outside the poly-n stretches of the masked sequence produced by the masking program.

The polynucleotides of the subject invention are isolated and obtained in
10 substantial purity, generally as other than an intact chromosome. Usually, the polynucleotides, either as DNA or RNA, will be obtained substantially free of other naturally-occurring nucleic acid sequences, generally being at least about 50%, usually at least about 90% pure and are typically "recombinant", *e.g.*, flanked by one or more nucleotides with which it is not normally associated on a naturally occurring
15 chromosome.

The polynucleotides of the invention can be provided as a linear molecule or within a circular molecule, and can be provided within autonomously replicating molecules (vectors) or within molecules without replication sequences. Expression of the polynucleotides can be regulated by their own or by other regulatory
20 sequences known in the art. The polynucleotides of the invention can be introduced into suitable host cells using a variety of techniques available in the art, such as transferrin polycation-mediated DNA transfer, transfection with naked or encapsulated nucleic acids, liposome-mediated DNA transfer, intracellular transportation of DNA-coated latex beads, protoplast fusion, viral infection, electroporation, gene gun, calcium
25 phosphate-mediated transfection, and the like.

The subject nucleic acid compositions can be used to, for example, produce polypeptides, as probes for the detection of mRNA of the invention in biological samples (*e.g.*, extracts of human cells) to generate additional copies of the polynucleotides, to generate ribozymes or antisense oligonucleotides, and as single
30 stranded DNA probes or as triple-strand forming oligonucleotides. The probes described herein can be used to, for example, determine the presence or absence of the polynucleotide sequences as shown in SEQ ID NOs:1-3351 or variants thereof in a sample. These and other uses are described in more detail below.

Use of Polynucleotides to Obtain Full-Length cDNA, Gene, and Promoter Region

Full-length cDNA molecules comprising the disclosed polynucleotides are obtained as follows. A polynucleotide having a sequence of one of SEQ ID NOs:1-3351, or a portion thereof comprising at least 12, 15, 18, or 20 nt, is used as a hybridization probe to detect hybridizing members of a cDNA library using probe design methods, cloning methods, and clone selection techniques such as those described in U.S. Patent No. 5,654,173. Libraries of cDNA are made from selected tissues, such as normal or tumor tissue, or from tissues of a mammal treated with, for example, a pharmaceutical agent. Preferably, the tissue is the same as the tissue from which the polynucleotides of the invention were isolated, as both the polynucleotides described herein and the cDNA represent expressed genes. Most preferably, the cDNA library is made from the biological material described herein in the Examples. The choice of cell type for library construction can be made after the identity of the protein encoded by the gene corresponding to the polynucleotide of the invention is known. This will indicate which tissue and cell types are likely to express the related gene, and thus represent a suitable source for the mRNA for generating the cDNA. As described in the Examples, cDNA of the invention was isolated from specific cell or tissue types, and such cells and tissues are preferable for obtaining related nucleic acids.

Techniques for producing and probing nucleic acid sequence libraries are described, for example, in Sambrook et al., *Molecular Cloning: A Laboratory Manual, 2nd Ed.*, (1989) Cold Spring Harbor Press, Cold Spring Harbor, NY. The cDNA can be prepared by using primers based on sequence from SEQ ID NOs:1-3351. In one embodiment, the cDNA library can be made from only poly-adenylated mRNA. Thus, poly-T primers can be used to prepare cDNA from the mRNA.

Members of the library that are larger than the provided polynucleotides, and preferably that encompass the complete coding sequence of the native message, are obtained. In order to confirm that the entire cDNA has been obtained, RNA protection experiments are performed as follows. Hybridization of a full-length cDNA to an mRNA will protect the RNA from RNase degradation. If the cDNA is not full length, then the portions of the mRNA that are not hybridized will be subject to RNase degradation. This is assayed, as is known in the art, by changes in electrophoretic mobility on polyacrylamide gels, or by detection of released monoribonucleotides. Sambrook et al., *Molecular Cloning: A Laboratory Manual, 2nd Ed.*, (1989) Cold Spring Harbor Press, Cold Spring Harbor, NY. In order to obtain additional sequences

5' to the end of a partial cDNA, 5' RACE (*PCR Protocols: A Guide to Methods and Applications*, (1990) Academic Press, Inc.) can be performed.

Genomic DNA is isolated using the provided polynucleotides in a manner similar to the isolation of full-length cDNAs. Briefly, the provided polynucleotides, or portions thereof, are used as probes to libraries of genomic DNA. Preferably, the library is obtained from the cell type that was used to generate the polynucleotides of the invention, but this is not essential. Most preferably, the genomic DNA is obtained from the biological material described herein in the Examples. Such libraries can be in vectors suitable for carrying large segments of a genome, such as P1 or YAC, as described in detail in Sambrook et al., 9.4-9.30. In addition, genomic sequences can be isolated from human BAC libraries, which are commercially available from Research Genetics, Inc., Huntsville, Alabama, USA, for example. In order to obtain additional 5' or 3' sequences, chromosome walking is performed, as described in Sambrook et al., such that adjacent and overlapping fragments of genomic DNA are isolated. These are mapped and pieced together, as is known in the art, using restriction digestion enzymes and DNA ligase.

Using the polynucleotide sequences of the invention, corresponding full-length genes can be isolated using both classical and PCR methods to construct and probe cDNA libraries. Using either method, Northern blots, preferably, are performed on a number of cell types to determine which cell lines express the gene of interest at the highest level. Classical methods of constructing cDNA libraries are taught in Sambrook et al., *supra*. With these methods, cDNA can be produced from mRNA and inserted into viral or expression vectors. Typically, libraries of mRNA comprising poly(A) tails can be produced with poly(T) primers. Similarly, cDNA libraries can be produced using the instant sequences as primers.

PCR methods are used to amplify the members of a cDNA library that comprise the desired insert. In this case, the desired insert will contain sequence from the full length cDNA that corresponds to the instant polynucleotides. Such PCR methods include gene trapping and RACE methods as described in Gruber et al., WO 95/04745 and Gruber et al., U.S. Patent No. 5,500,356. Kits are commercially available to perform gene trapping experiments from, for example, Life Technologies, Gaithersburg, Maryland, USA. In preferred embodiments of RACE, a common primer is designed to anneal to an arbitrary adaptor sequence ligated to cDNA ends (Apte and Siebert, *Biotechniques* (1993) 15:890-893; Edwards et al., *Nuc. Acids Res.* (1991) 19:5227-5232). When a single gene-specific RACE primer is paired with the common

primer, preferential amplification of sequences between the single gene specific primer and the common primer occurs. Commercial cDNA pools modified for use in RACE are available.

The promoter region of a gene generally is located 5' to the initiation site for RNA polymerase II. Hundreds of promoter regions contain the "TATA" box, a sequence such as TATTA or TATAA, which is sensitive to mutations. The promoter region can be obtained by performing 5' RACE using a primer from the coding region of the gene. Alternatively, the cDNA can be used as a probe for the genomic sequence, and the region 5' to the coding region is identified by "walking up." If the gene is highly expressed or differentially expressed, the promoter from the gene can be of use in a regulatory construct for a heterologous gene.

Once the full-length cDNA or gene is obtained, DNA encoding variants can be prepared by site-directed mutagenesis, described in detail in Sambrook et al., 15.3-15.63. The choice of codon or nucleotide to be replaced can be based on disclosure herein on optional changes in amino acids to achieve altered protein structure and/or function.

As an alternative method to obtaining DNA or RNA from a biological material, nucleic acid comprising nucleotides having the sequence of one or more polynucleotides of the invention can be synthesized. Thus, the invention encompasses nucleic acid molecules ranging in length from 15 nt (corresponding to at least 15 contiguous nt of one of SEQ ID NOs:1-3351) up to a maximum length suitable for one or more biological manipulations, including replication and expression, of the nucleic acid molecule. The invention includes but is not limited to (a) nucleic acid having the size of a full gene, and comprising at least one of SEQ ID NOs:1-3351; (b) the nucleic acid of (a) also comprising at least one additional polynucleotide or gene, operably linked to permit expression of a fusion protein; (c) an expression vector comprising (a) or (b); (d) a plasmid comprising (a) or (b); and (e) a recombinant viral particle comprising (a) or (b). Once provided with the polynucleotides disclosed herein, construction or preparation of (a) - (e) are well within the skill in the art.

The sequence of a nucleic acid comprising at least 15 contiguous nt of at least any one of SEQ ID NOs:1-3351, preferably the entire sequence of at least any one of SEQ ID NOs:1-3351, is not limited and can be any sequence of A, T, G, and/or C (for DNA) and A, U, G, and/or C (for RNA) or modified bases thereof, including inosine and pseudouridine. The choice of sequence will depend on the desired function and can be dictated by coding regions desired, the intron-like regions desired, and the

regulatory regions desired. Where the entire sequence of any one of SEQ ID NOs:1-3351 is within the nucleic acid, the nucleic acid obtained is referred to herein as a polynucleotide comprising the sequence of any one of SEQ ID NOs:1-3351.

Expression of Polypeptide Encoded by Full-Length cDNA or Full-Length Gene

5 The provided polynucleotides (e.g., a polynucleotide having a sequence of one of SEQ ID NOs:1-3351), the corresponding cDNA, or the full-length gene is used to express a partial or complete gene product. Constructs of polynucleotides having sequences of SEQ ID NOs:1-3351 can be generated synthetically. Alternatively, single-step assembly of a gene and entire plasmid from large numbers of
10 oligodeoxyribonucleotides is described by, e.g., Stemmer et al., *Gene (Amsterdam)* (1995) 164(1):49-53. In this method, assembly PCR (the synthesis of long DNA sequences from large numbers of oligodeoxyribonucleotides (oligos)) is described. The method is derived from DNA shuffling (Stemmer, *Nature* (1994) 370:389-391), and does not rely on DNA ligase, but instead relies on DNA polymerase to build
15 increasingly longer DNA fragments during the assembly process.

Appropriate polynucleotide constructs are purified using standard recombinant DNA techniques as described in, for example, Sambrook et al., *Molecular Cloning: A Laboratory Manual, 2nd Ed.*, (1989) Cold Spring Harbor Press, Cold Spring Harbor, NY, and under current regulations described in United States Dept. of HHS,
20 National Institute of Health (NIH) Guidelines for Recombinant DNA Research. The gene product encoded by a polynucleotide of the invention is expressed in any expression system, including, for example, bacterial, yeast, insect, amphibian and mammalian systems. Vectors, host cells and methods for obtaining expression in same are well known in the art. Suitable vectors and host cells are described in U.S. Patent
25 No. 5,654,173.

Polynucleotide molecules comprising a polynucleotide sequence provided herein are generally propagated by placing the molecule in a vector. Viral and non-viral vectors are used, including plasmids. The choice of plasmid will depend on the type of cell in which propagation is desired and the purpose of propagation. Certain
30 vectors are useful for amplifying and making large amounts of the desired DNA sequence. Other vectors are suitable for expression in cells in culture. Still other vectors are suitable for transfer and expression in cells in a whole animal or person. The choice of appropriate vector is well within the skill of the art. Many such vectors are

available commercially. Methods for preparation of vectors comprising a desired sequence are well known in the art.

The polynucleotides set forth in SEQ ID NOs:1-3351 or their corresponding full-length polynucleotides are linked to regulatory sequences as appropriate to obtain the desired expression properties. These can include promoters (attached either at the 5' end of the sense strand or at the 3' end of the antisense strand), enhancers, terminators, operators, repressors, and inducers. The promoters can be regulated or constitutive. In some situations it may be desirable to use conditionally active promoters, such as tissue-specific or developmental stage-specific promoters. These are linked to the desired nucleotide sequence using the techniques described above for linkage to vectors. Any techniques known in the art can be used.

When any appropriate host cells or organisms are used to replicate and/or express the polynucleotides or nucleic acids of the invention, the resulting replicated nucleic acid, RNA, expressed protein or polypeptide, is within the scope of the invention as a product of the host cell or organism. The product is recovered by any appropriate means known in the art.

Once the gene corresponding to a selected polynucleotide is identified, its expression can be regulated in the cell to which the gene is native. For example, an endogenous gene of a cell can be regulated by an exogenous regulatory sequence as disclosed in U.S. Patent No. 5,641,670.

Identification of Functional and Structural Motifs of Novel Genes

Translations of the nucleotide sequence of the provided polynucleotides, cDNAs or full genes can be aligned with individual known sequences. Similarity with individual sequences can be used to determine the activity of the polypeptides encoded by the polynucleotides of the invention. Also, sequences exhibiting similarity with more than one individual sequence can exhibit activities that are characteristic of either or both individual sequences.

The full length sequences and fragments of the polynucleotide sequences of the nearest neighbors can be used as probes and primers to identify and isolate the full length sequence corresponding to provided polynucleotides. The nearest neighbors can indicate a tissue or cell type to be used to construct a library for the full-length sequences corresponding to the provided polynucleotides.

Typically, a selected polynucleotide is translated in all six frames to determine the best alignment with the individual sequences. The sequences disclosed

herein in the Sequence Listing are in a 5' to 3' orientation and translation in three frames can be sufficient. These amino acid sequences are referred to, generally, as query sequences, which will be aligned with the individual sequences. Databases with individual sequences are described in "Computer Methods for Macromolecular
5 Sequence Analysis" *Methods in Enzymology* (1996) 266, Doolittle, Academic Press, Inc., a division of Harcourt Brace & Co., San Diego, California, USA. Databases include Genbank, EMBL, and DNA Database of Japan (DDBJ).

Query and individual sequences can be aligned using the methods and computer programs described above, and include BLAST, available over the world
10 wide web at <http://www.ncbi.nlm.nih.gov/BLAST>. Another alignment algorithm is Fasta, available in the Genetics Computing Group (GCG) package, Madison, Wisconsin, USA, a wholly owned subsidiary of Oxford Molecular Group, Inc. Other techniques for alignment are described in Doolittle, *supra*. Preferably, an alignment program that permits gaps in the sequence is utilized to align the sequences. The
15 Smith-Waterman is one type of algorithm that permits gaps in sequence alignments. See *Meth. Mol. Biol.* (1997) 70: 173-187. Also, the GAP program using the Needleman and Wunsch alignment method can be utilized to align sequences. An alternative search strategy uses MPSRCH software, which runs on a MASPAR computer. MPSRCH uses a Smith-Waterman algorithm to score sequences on a massively parallel computer.
20 This approach improves ability to identify sequences that are distantly related matches, and is especially tolerant of small gaps and nucleotide sequence errors. Amino acid sequences encoded by the provided polynucleotides can be used to search both protein and DNA databases.

High Similarity. In general, in alignment results considered to be of high
25 similarity, the percent of the alignment region length is typically at least about 55% of total length query sequence; more typically, at least about 58%; even more typically; at least about 60% of the total residue length of the query sequence. Usually, percent length of the alignment region can be as much as about 62%; more usually, as much as about 64%; even more usually, as much as about 66%. Further, for high similarity, the
30 region of alignment, typically, exhibits at least about 75% of sequence identity; more typically, at least about 78%; even more typically; at least about 80% sequence identity. Usually, percent sequence identity can be as much as about 82%; more usually, as much as about 84%; even more usually, as much as about 86%.

The p value is used in conjunction with these methods. If high similarity
35 is found, the query sequence is considered to have high similarity with a profile

sequence when the p value is less than or equal to about 10^{-2} ; more usually; less than or equal to about 10^{-3} ; even more usually; less than or equal to about 10^{-4} . More typically, the p value is no more than about 10^{-5} ; more typically; no more than or equal to about 10^{-10} ; even more typically; no more than or equal to about 10^{-15} for the query sequence
5 to be considered high similarity.

Similarity Determined by Sequence Identity Alone. Sequence identity alone can be used to determine similarity of a query sequence to an individual sequence and can indicate the activity of the sequence. Such an alignment, preferably, permits gaps to align sequences. Typically, the query sequence is related to the profile sequence
10 if the sequence identity over the entire query sequence is at least about 15%; more typically, at least about 20%; even more typically, at least about 25%; even more typically, at least about 50%. Sequence identity alone as a measure of similarity is most useful when the query sequence is usually, at least 80 residues in length; more usually, 90 residues; even more usually, at least 95 amino acid residues in length. More
15 typically, similarity can be concluded based on sequence identity alone when the query sequence is preferably 100 residues in length; more preferably, 120 residues in length; even more preferably, 150 amino acid residues in length.

Alignments with Profile and Multiple Aligned Sequences. Translations of the provided polynucleotides can be aligned with amino acid profiles that define
20 either protein families or common motifs. Also, translations of the provided polynucleotides can be aligned to multiple sequence alignments (MSA) comprising the polypeptide sequences of members of protein families or motifs. Similarity or identity with profile sequences or MSAs can be used to determine the activity of the gene products (e.g., polypeptides) encoded by the provided polynucleotides or corresponding
25 cDNA or genes. For example, sequences that show an identity or similarity with a chemokine profile or MSA can exhibit chemokine activities.

Profiles can be designed manually by (1) creating an MSA, which is an alignment of the amino acid sequence of members that belong to the family and (2) constructing a statistical representation of the alignment. Such methods are described,
30 for example, in Birney et al., *Nucl. Acid Res.* (1996) 24(14): 2730-2739. MSAs of some protein families and motifs are publicly available. MSAs are described also in Sonnhammer et al., *Proteins* (1997) 28: 405-420. A brief description of MSAs is reported in Pascarella et al., *Prot. Eng.* (1996) 9(3):249-251. Techniques for building profiles from MSAs are described in Sonnhammer et al., *supra*; Birney et al., *supra*;

and "Computer Methods for Macromolecular Sequence Analysis," *Methods in Enzymology* (1996) 266, Doolittle, Academic Press, Inc., San Diego, California, USA.

Similarity between a query sequence and a protein family or motif can be determined by (a) comparing the query sequence against the profile and/or (b) aligning the query sequence with the members of the family or motif. Typically, a program such as Searchwise is used to compare the query sequence to the statistical representation of the multiple alignment, also known as a profile (see Birney et al., *supra*). Other techniques to compare the sequence and profile are described in Sonnhammer et al., *supra* and Doolittle, *supra*.

Next, methods described by Feng et al., *J. Mol. Evol.* (1987) 25:351 and Higgins et al., *CABIOS* (1989) 5:151 can be used align the query sequence with the members of a family or motif, also known as a MSA. Sequence alignments can be generated using any of a variety of software tools. Examples include PileUp, which creates a multiple sequence alignment, and is described in Feng et al., *J. Mol. Evol.* (1987) 25:351. Another method, GAP, uses the alignment method of Needleman et al., *J. Mol. Biol.* (1970) 48:443. GAP is best suited for global alignment of sequences. A third method, BestFit, functions by inserting gaps to maximize the number of matches using the local homology algorithm of Smith et al., *Adv. Appl. Math.* (1981) 2:482. In general, the following factors are used to determine if a similarity between a query sequence and a profile or MSA exists: (1) number of conserved residues found in the query sequence, (2) percentage of conserved residues found in the query sequence, (3) number of frameshifts, and (4) spacing between conserved residues.

Some alignment programs that both translate and align sequences can make any number of frameshifts when translating the nucleotide sequence to produce the best alignment. The fewer frameshifts needed to produce an alignment, the stronger the similarity or identity between the query and profile or MSAs. For example, a weak similarity resulting from no frameshifts can be a better indication of activity or structure of a query sequence, than a strong similarity resulting from two frameshifts. Preferably, three or fewer frameshifts are found in an alignment; more preferably two or fewer frameshifts; even more preferably, one or fewer frameshifts; even more preferably, no frameshifts are found in an alignment of query and profile or MSAs.

Conserved residues are those amino acids found at a particular position in all or some of the family or motif members. Alternatively, a position is considered conserved if only a certain class of amino acids is found in a particular position in all or

some of the family members. For example, the N-terminal position can contain a positively charged amino acid, such as lysine, arginine, or histidine.

Typically, a residue of a polypeptide is conserved when a class of amino acids or a single amino acid is found at a particular position in at least about 40% of all class members; more typically, at least about 50%; even more typically, at least about 60% of the members. Usually, a residue is conserved when a class or single amino acid is found in at least about 70% of the members of a family or motif; more usually, at least about 80%; even more usually, at least about 90%; even more usually, at least about 95%.

A residue is considered conserved when three unrelated amino acids are found at a particular position in the some or all of the members; more usually, two unrelated amino acids. These residues are conserved when the unrelated amino acids are found at particular positions in at least about 40% of all class member; more typically, at least about 50%; even more typically, at least about 60% of the members. Usually, a residue is conserved when a class or single amino acid is found in at least about 70% of the members of a family or motif; more usually, at least about 80%; even more usually, at least about 90%; even more usually, at least about 95%.

A query sequence has similarity to a profile or MSA when the query sequence comprises at least about 25% of the conserved residues of the profile or MSA; more usually, at least about 30%; even more usually, at least about 40%. Typically, the query sequence has a stronger similarity to a profile sequence or MSA when the query sequence comprises at least about 45% of the conserved residues of the profile or MSA; more typically, at least about 50%; even more typically, at least about 55%.

Identification of Secreted and Membrane-Bound Polypeptides

Both secreted and membrane-bound polypeptides of the present invention are of particular interest. For example, levels of secreted polypeptides can be assayed in body fluids that are convenient, such as blood, plasma, serum, and other body fluids such as urine, prostatic fluid and semen. Membrane-bound polypeptides are useful for constructing vaccine antigens or inducing an immune response. Such antigens would comprise all or part of the extracellular region of the membrane-bound polypeptides. Because both secreted and membrane-bound polypeptides comprise a fragment of contiguous hydrophobic amino acids, hydrophobicity predicting algorithms can be used to identify such polypeptides.

A signal sequence is usually encoded by both secreted and membrane-bound polypeptide genes to direct a polypeptide to the surface of the cell. The signal sequence usually comprises a stretch of hydrophobic residues. Such signal sequences can fold into helical structures. Membrane-bound polypeptides typically comprise at least one transmembrane region that possesses a stretch of hydrophobic amino acids that can transverse the membrane. Some transmembrane regions also exhibit a helical structure. Hydrophobic fragments within a polypeptide can be identified by using computer algorithms. Such algorithms include Hopp & Woods, *Proc. Natl. Acad. Sci. USA* (1981) 78:3824-3828; Kyte & Doolittle, *J. Mol. Biol.* (1982) 157: 105-132; and RAOAR algorithm, Degli Esposti et al., *Eur. J. Biochem.* (1990) 190: 207-219.

Another method of identifying secreted and membrane-bound polypeptides is to translate the polynucleotides of the invention in all six frames and determine if at least 8 contiguous hydrophobic amino acids are present. Those translated polypeptides with at least 8; more typically, 10; even more typically, 12 contiguous hydrophobic amino acids are considered to be either a putative secreted or membrane bound polypeptide. Hydrophobic amino acids include alanine, glycine, histidine, isoleucine, leucine, lysine, methionine, phenylalanine, proline, threonine, tryptophan, tyrosine, and valine

Identification of the Function of an Expression Product of a Full-Length Gene

Ribozymes, antisense constructs, and dominant negative mutants can be used to determine function of the expression product of a gene corresponding to a polynucleotide provided herein. The phosphoramidite method of oligonucleotide synthesis can be used to construct antisense molecules and ribozymes. See Beaucage et al., *Tet. Lett.* (1981) 22:1859 and U.S. Patent No. 4,668,777. Automated devices for synthesis are available to create oligonucleotides using this chemistry. Examples of such devices include Biosearch 8600, Models 392 and 394 by Applied Biosystems, a division of Perkin-Elmer Corp., Foster City, California, USA; and Expedite by Perceptive Biosystems, Framingham, Massachusetts, USA. Synthetic RNA, phosphate analog oligonucleotides, and chemically derivatized oligonucleotides can also be produced, and can be covalently attached to other molecules. RNA oligonucleotides can be synthesized, for example, using RNA phosphoramidites. This method can be performed on an automated synthesizer, such as Applied Biosystems, Models 392 and 394, Foster City, California, USA.

18

Oligonucleotides of up to 200 nt can be synthesized, more typically, 100 nt, more typically 50 nt; even more typically 30 to 40 nt. These synthetic fragments can be annealed and ligated together to construct larger fragments. See, for example, Sambrook et al., *supra*. Trans-cleaving catalytic RNAs (ribozymes) are RNA molecules possessing endoribonuclease activity. Ribozymes are specifically designed for a particular target, and the target message must contain a specific nucleotide sequence. They are engineered to cleave any RNA species site-specifically in the background of cellular RNA. The cleavage event renders the mRNA unstable and prevents protein expression. Importantly, ribozymes can be used to inhibit expression of a gene of unknown function for the purpose of determining its function in an *in vitro* or *in vivo* context, by detecting the phenotypic effect.

Antisense nucleic acids are designed to specifically bind to RNA, resulting in the formation of RNA-DNA or RNA-RNA hybrids, with an arrest of DNA replication, reverse transcription or messenger RNA translation. Antisense polynucleotides based on a selected polynucleotide sequence can interfere with expression of the corresponding gene. Antisense polynucleotides are typically generated within the cell by expression from antisense constructs that contain the antisense strand as the transcribed strand. Antisense polynucleotides based on the disclosed polynucleotides will bind and/or interfere with the translation of mRNA comprising a sequence complementary to the antisense polynucleotide. The expression products of control cells and cells treated with the antisense construct are compared to detect the protein product of the gene corresponding to the polynucleotide upon which the antisense construct is based. The protein is isolated and identified using routine biochemical methods.

Given the extensive background literature and clinical experience in antisense therapy, one skilled in the art can use selected polynucleotides of the invention as additional potential therapeutics. The choice of polynucleotide can be narrowed by first testing them for binding to "hot spot" regions of the genome of cancerous cells. If a polynucleotide is identified as binding to a "hot spot," testing the polynucleotide as an antisense compound in the corresponding cancer cells is warranted.

Dominant negative mutations also are readily generated for corresponding proteins that are active as homomultimers. A mutant polypeptide will interact with wild-type polypeptides (made from the other allele) and form a non-functional multimer. Thus, a mutation is in a substrate-binding domain, a catalytic

domain, or a cellular localization domain. Preferably, the mutant polypeptide will be overproduced. Point mutations are made that have such an effect. In addition, fusion of different polypeptides of various lengths to the terminus of a protein can yield dominant negative mutants. General strategies are available for making dominant negative
5 mutants (see, e.g., Herskowitz, *Nature* (1987) 329:219). Such techniques can be used to create loss of function mutations, which are useful for determining protein function.

Polypeptides and Variants Thereof

The polypeptides of the invention include those encoded by the disclosed polynucleotides, as well as nucleic acids that, by virtue of the degeneracy of the genetic
10 code, are not identical in sequence to the disclosed polynucleotides. Thus, the invention includes within its scope a polypeptide encoded by a polynucleotide having the sequence of any one of SEQ ID NOs:1-3351 or a variant thereof.

In general, the term "polypeptide" as used herein refers to both the full length polypeptide encoded by the recited polynucleotide, the polypeptide encoded by
15 the gene represented by the recited polynucleotide, as well as portions or fragments thereof. "Polypeptides" also includes variants of the naturally occurring proteins, where such variants are homologous or substantially similar to the naturally occurring protein, and can be of an origin of the same or different species as the naturally occurring protein (e.g., human, murine, or some other species that naturally expresses the recited
20 polypeptide, usually a mammalian species). In general, variant polypeptides have a sequence that has at least about 80%, usually at least about 90%, and more usually at least about 98% sequence identity with a differentially expressed polypeptide of the invention, as measured by BLAST using the parameters described above. The variant polypeptides can be naturally or non-naturally glycosylated, i.e., the polypeptide has a
25 glycosylation pattern that differs from the glycosylation pattern found in the corresponding naturally occurring protein.

The invention also encompasses homologs of the disclosed polypeptides (or fragments thereof) where the homologs are isolated from other species, i.e., other animal or plant species, where such homologs, usually mammalian species, e.g.,
30 rodents, such as mice, rats; domestic animals, e.g., horse, cow, dog, cat; and humans. By "homolog" is meant a polypeptide having at least about 35%, usually at least about 40% and more usually at least about 60% amino acid sequence identity to a particular differentially expressed protein as identified above, where sequence identity is determined using the BLAST algorithm, with the parameters described above.

In general, the polypeptides of the subject invention are provided in a non-naturally occurring environment, e.g., are separated from their naturally occurring environment. In certain embodiments, the subject protein is present in a composition that is enriched for the protein as compared to a control. As such, purified polypeptide
5 is provided, where by purified is meant that the protein is present in a composition that is substantially free of non-differentially expressed polypeptides, where by substantially free is meant that less than 90%, usually less than 60% and more usually less than 50% of the composition is made up of non-differentially expressed polypeptides.

Also within the scope of the invention are variants; variants of
10 polypeptides include mutants, fragments, and fusions. Mutants can include amino acid substitutions, additions or deletions. The amino acid substitutions can be conservative amino acid substitutions or substitutions to eliminate non-essential amino acids, such as to alter a glycosylation site, a phosphorylation site or an acetylation site, or to minimize misfolding by substitution or deletion of one or more cysteine residues that are not
15 necessary for function. Conservative amino acid substitutions are those that preserve the general charge, hydrophobicity/ hydrophilicity, and/or steric bulk of the amino acid substituted. Variants can be designed so as to retain biological activity of a particular region of the protein (e.g., a functional domain and/or, where the polypeptide is a member of a protein family, a region associated with a consensus sequence). Selection
20 of amino acid alterations for production of variants can be based upon the accessibility (interior vs. exterior) of the amino acid (see, e.g., Go et al., *Int. J. Peptide Protein Res.* (1980) 15:211), the thermostability of the variant polypeptide (see, e.g., Querol et al., *Prot. Eng.* (1996) 9:265), desired glycosylation sites (see, e.g., Olsen and Thomsen, *J. Gen. Microbiol.* (1991) 137:579), desired disulfide bridges (see, e.g., Clarke et al.,
25 *Biochemistry* (1993) 32:4322; and Wakarchuk et al., *Protein Eng.* (1994) 7:1379), desired metal binding sites (see, e.g., Toma et al., *Biochemistry* (1991) 30:97, and Haezrebrouck et al., *Protein Eng.* (1993) 6:643), and desired substitutions with in proline loops (see, e.g., Masul et al., *Appl. Env. Microbiol.* (1994) 60:3579). Cysteine-depleted muteins can be produced as disclosed in U.S. Patent No. 4,959,314.

30 Variants also include fragments of the polypeptides disclosed herein, particularly biologically active fragments and/or fragments corresponding to functional domains. Fragments of interest will typically be at least about 10 aa to at least about 15 aa in length, usually at least about 50 aa in length, and can be as long as 300 aa in length or longer, but will usually not exceed about 1000 aa in length, where the fragment will
35 have a stretch of amino acids that is identical to a polypeptide encoded by a

21

polynucleotide having a sequence of any SEQ ID NOs:1-3351, or a homolog thereof. The protein variants described herein are encoded by polynucleotides that are within the scope of the invention. The genetic code can be used to select the appropriate codons to construct the corresponding variants.

5 Computer-Related Embodiments

 In general, a library of polynucleotides is a collection of sequence information, which information is provided in either biochemical form (*e.g.*, as a collection of polynucleotide molecules), or in electronic form (*e.g.*, as a collection of polynucleotide sequences stored in a computer-readable form, as in a computer system and/or as part of a computer program). The sequence information of the polynucleotides can be used in a variety of ways, *e.g.*, as a resource for gene discovery, as a representation of sequences expressed in a selected cell type (*e.g.*, cell type markers), and/or as markers of a given disease or disease state. In general, a disease marker is a representation of a gene product that is present in all cells affected by disease either at an increased or decreased level relative to a normal cell (*e.g.*, a cell of the same or similar type that is not substantially affected by disease). For example, a polynucleotide sequence in a library can be a polynucleotide that represents an mRNA, polypeptide, or other gene product encoded by the polynucleotide, that is either overexpressed or underexpressed in a breast ductal cell affected by cancer relative to a normal (*i.e.*, substantially disease-free) breast cell.

 The nucleotide sequence information of the library can be embodied in any suitable form, *e.g.*, electronic or biochemical forms. For example, a library of sequence information embodied in electronic form comprises an accessible computer data file (or, in biochemical form, a collection of nucleic acid molecules) that contains the representative nucleotide sequences of genes that are differentially expressed (*e.g.*, overexpressed or underexpressed) as between, for example, i) a cancerous cell and a normal cell; ii) a cancerous cell and a dysplastic cell; iii) a cancerous cell and a cell affected by a disease or condition other than cancer; iv) a metastatic cancerous cell and a normal cell and/or non-metastatic cancerous cell; v) a malignant cancerous cell and a non-malignant cancerous cell (or a normal cell) and/or vi) a dysplastic cell relative to a normal cell. Other combinations and comparisons of cells affected by various diseases or stages of disease will be readily apparent to the ordinarily skilled artisan. Biochemical embodiments of the library include a collection of nucleic acids that have

the sequences of the genes in the library, where the nucleic acids can correspond to the entire gene in the library or to a fragment thereof, as described in greater detail below.

The polynucleotide libraries of the subject invention generally comprise sequence information of a plurality of polynucleotide sequences, where at least one of the polynucleotides has a sequence of any of SEQ ID NOs:1-3351. By plurality is meant at least 2, usually at least 3 and can include up to all of SEQ ID NOs:1-3351. The length and number of polynucleotides in the library will vary with the nature of the library, *e.g.*, if the library is an oligonucleotide array, a cDNA array, a computer database of the sequence information, *etc.*

Where the library is an electronic library, the nucleic acid sequence information can be present in a variety of media. "Media" refers to a manufacture, other than an isolated nucleic acid molecule, that contains the sequence information of the present invention. Such a manufacture provides the genome sequence or a subset thereof in a form that can be examined by means not directly applicable to the sequence as it exists in a nucleic acid. For example, the nucleotide sequence of the present invention, *e.g.*, the nucleic acid sequences of any of the polynucleotides of SEQ ID NOs:1-3351, can be recorded on computer readable media, *e.g.*, any medium that can be read and accessed directly by a computer. Such media include, but are not limited to: magnetic storage media, such as a floppy disc, a hard disc storage medium, and a magnetic tape; optical storage media such as CD-ROM; electrical storage media such as RAM and ROM; and hybrids of these categories such as magnetic/optical storage media. One of skill in the art can readily appreciate how any of the presently known computer readable mediums can be used to create a manufacture comprising a recording of the present sequence information. "Recorded" refers to a process for storing information on computer readable medium, using any such methods as known in the art. Any convenient data storage structure can be chosen, based on the means used to access the stored information. A variety of data processor programs and formats can be used for storage, *e.g.*, word processing text file, database format, *etc.* In addition to the sequence information, electronic versions of the libraries of the invention can be provided in conjunction or connection with other computer-readable information and/or other types of computer-readable files (*e.g.*, searchable files, executable files, *etc.*, including, but not limited to, for example, search program software, *etc.*).

By providing the nucleotide sequence in computer readable form, the information can be accessed for a variety of purposes. Computer software to access sequence information is publicly available. For example, the BLAST (Altschul et al.,

supra.) and BLAZE (Brutlag et al. *Comp. Chem.* (1993) 17:203) search algorithms on a Sybase system can be used to identify open reading frames (ORFs) within the genome that contain homology to ORFs from other organisms.

As used herein, "a computer-based system" refers to the hardware means, software means, and data storage means used to analyze the nucleotide sequence information of the present invention. The minimum hardware of the computer-based systems of the present invention comprises a central processing unit (CPU), input means, output means, and data storage means. A skilled artisan can readily appreciate that any one of the currently available computer-based system are suitable for use in the present invention. The data storage means can comprise any manufacture comprising a recording of the present sequence information as described above, or a memory access means that can access such a manufacture.

"Search means" refers to one or more programs implemented on the computer-based system, to compare a target sequence or target structural motif, or expression levels of a polynucleotide in a sample, with the stored sequence information. Search means can be used to identify fragments or regions of the genome that match a particular target sequence or target motif. A variety of known algorithms are publicly known and commercially available, *e.g.*, MacPattern (EMBL), BLASTN and BLASTX (NCBI). A "target sequence" can be any polynucleotide or amino acid sequence of six or more contiguous nucleotides or two or more amino acids, preferably from about 10 to 100 amino acids or from about 30 to 300 nt. A variety of comparing means can be used to accomplish comparison of sequence information from a sample (*e.g.*, to analyze target sequences, target motifs, or relative expression levels) with the data storage means. A skilled artisan can readily recognize that any one of the publicly available homology search programs can be used as the search means for the computer based systems of the present invention to accomplish comparison of target sequences and motifs. Computer programs to analyze expression levels in a sample and in controls are also known in the art.

A "target structural motif," or "target motif," refers to any rationally selected sequence or combination of sequences in which the sequence(s) are chosen based on a three-dimensional configuration that is formed upon the folding of the target motif, or on consensus sequences of regulatory or active sites. There are a variety of target motifs known in the art. Protein target motifs include, but are not limited to, enzyme active sites and signal sequences. Nucleic acid target motifs include, but are

not limited to, hairpin structures, promoter sequences and other expression elements such as binding sites for transcription factors.

5 A variety of structural formats for the input and output means can be used to input and output the information in the computer-based systems of the present invention. One format for an output means ranks the relative expression levels of different polynucleotides. Such presentation provides a skilled artisan with a ranking of relative expression levels to determine a gene expression profile.

10 As discussed above, the "library" of the invention also encompasses biochemical libraries of the polynucleotides of SEQ ID NOs:1-3351, *e.g.*, collections of nucleic acids representing the provided polynucleotides. The biochemical libraries can take a variety of forms, *e.g.*, a solution of cDNAs, a pattern of probe nucleic acids stably associated with a surface of a solid support (*i.e.*, an array) and the like. Of particular interest are nucleic acid arrays in which one or more of SEQ ID NOs:1-3351 is represented on the array. By array is meant an article of manufacture that has at least a
15 substrate with at least two distinct nucleic acid targets on one of its surfaces, where the number of distinct nucleic acids can be considerably higher, typically being at least 10 nt, usually at least 20 nt and often at least 25 nt. A variety of different array formats have been developed and are known to those of skill in the art. The arrays of the subject invention find use in a variety of applications, including gene expression analysis, drug
20 screening, mutation analysis and the like, as disclosed in the above-listed exemplary patent documents.

In addition to the above nucleic acid libraries, analogous libraries of polypeptides are also provided, where the where the polypeptides of the library will represent at least a portion of the polypeptides encoded by SEQ ID NOs:1-3351.

25 Use of Polynucleotide Probes in Mapping, and in Tissue Profiling

Polynucleotide probes, generally comprising at least 12 contiguous nt of a polynucleotide as shown in the Sequence Listing, are used for a variety of purposes, such as chromosome mapping of the polynucleotide and detection of transcription levels. Additional disclosure about preferred regions of the disclosed polynucleotide
30 sequences is found in the Examples. A probe that hybridizes specifically to a polynucleotide disclosed herein should provide a detection signal at least 5-, 10-, or 20-fold higher than the background hybridization provided with other unrelated sequences.

Detection of Expression Levels. Nucleotide probes are used to detect expression of a gene corresponding to the provided polynucleotide. In Northern blots,

mRNA is separated electrophoretically and contacted with a probe. A probe is detected as hybridizing to an mRNA species of a particular size. The amount of hybridization is quantitated to determine relative amounts of expression, for example under a particular condition. Probes are used for *in situ* hybridization to cells to detect expression. Probes
5 can also be used *in vivo* for diagnostic detection of hybridizing sequences. Probes are typically labeled with a radioactive isotope. Other types of detectable labels can be used such as chromophores, fluors, and enzymes. Other examples of nucleotide hybridization assays are described in WO92/02526 and U.S. Patent No. 5,124,246.

Alternatively, the Polymerase Chain Reaction (PCR) is another means
10 for detecting small amounts of target nucleic acids (see, *e.g.*, Mullis et al., *Meth. Enzymol.* (1987) 155:335; U.S. Patent No. 4,683,195; and U.S. Patent No. 4,683,202). Two primer polynucleotides nucleotides that hybridize with the target nucleic acids are used to prime the reaction. The primers can be composed of sequence within or 3' and 5' to the polynucleotides of the Sequence Listing. Alternatively, if the primers are 3' and
15 5' to these polynucleotides, they need not hybridize to them or the complements. After amplification of the target with a thermostable polymerase, the amplified target nucleic acids can be detected by methods known in the art, *e.g.*, Southern blot. mRNA or cDNA can also be detected by traditional blotting techniques (*e.g.*, Southern blot, Northern blot, *etc.*) described in Sambrook et al., "Molecular Cloning: A Laboratory
20 Manual" (New York, Cold Spring Harbor Laboratory, 1989) (*e.g.*, without PCR amplification). In general, mRNA or cDNA generated from mRNA using a polymerase enzyme can be purified and separated using gel electrophoresis, and transferred to a solid support, such as nitrocellulose. The solid support is exposed to a labeled probe, washed to remove any unhybridized probe, and duplexes containing the labeled probe
25 are detected.

Mapping. Polynucleotides of the present invention can be used to identify a chromosome on which the corresponding gene resides. Such mapping can be useful in identifying the function of the polynucleotide-related gene by its proximity to other genes with known function. Function can also be assigned to the polynucleotide-
30 related gene when particular syndromes or diseases map to the same chromosome. For example, use of polynucleotide probes in identification and quantification of nucleic acid sequence aberrations is described in U.S. Patent No. 5,783,387. An exemplary mapping method is fluorescence *in situ* hybridization (FISH), which facilitates comparative genomic hybridization to allow total genome assessment of changes in
35 relative copy number of DNA sequences (see, *e.g.*, Valdes et al., *Methods in Molecular*

Biology (1997) 68:1). Polynucleotides can also be mapped to particular chromosomes using, for example, radiation hybrids or chromosome-specific hybrid panels. See Leach et al., *Advances in Genetics*, (1995) 33:63-99; Walter et al., *Nature Genetics* (1994) 7:22; Walter and Goodfellow, *Trends in Genetics* (1992) 9:352. Panels for radiation
5 hybrid mapping are available from Research Genetics, Inc., Huntsville, Alabama, USA. The statistical program RHMAP can be used to construct a map based on the data from radiation hybridization with a measure of the relative likelihood of one order versus another. RHMAP is available via the world wide web at <http://www.sph.umich.edu/group/statgen/software>. In addition, commercial programs are available for identifying
10 regions of chromosomes commonly associated with disease, such as cancer.

Tissue Typing or Profiling. Expression of specific mRNA corresponding to the provided polynucleotides can vary in different cell types and can be tissue-specific. This variation of mRNA levels in different cell types can be exploited with nucleic acid probe assays to determine tissue types. For example, PCR,
15 branched DNA probe assays, or blotting techniques utilizing nucleic acid probes substantially identical or complementary to polynucleotides listed in the Sequence Listing can determine the presence or absence of the corresponding cDNA or mRNA.

Tissue typing can be used to identify the developmental organ or tissue source of a metastatic lesion by identifying the expression of a particular marker of that
20 organ or tissue. If a polynucleotide is expressed only in a specific tissue type, and a metastatic lesion is found to express that polynucleotide, then the developmental source of the lesion has been identified. Expression of a particular polynucleotide can be assayed by detection of either the corresponding mRNA or the protein product.

Use of Polymorphisms. A polynucleotide of the invention can be used in
25 forensics, genetic analysis, mapping, and diagnostic applications where the corresponding region of a gene is polymorphic in the human population. Any means for detecting a polymorphism in a gene can be used, including, but not limited to electrophoresis of protein polymorphic variants, differential sensitivity to restriction enzyme cleavage, and hybridization to allele-specific probes.

30 Antibody Production

Expression products of a polynucleotide of the invention, as well as the corresponding mRNA, cDNA, or complete gene, can be prepared and used for raising antibodies for experimental, diagnostic, and therapeutic purposes. For polynucleotides to which a corresponding gene has not been assigned, this provides an additional

method of identifying the corresponding gene. The polynucleotide or related cDNA is expressed as described above, and antibodies are prepared. These antibodies are specific to an epitope on the polypeptide encoded by the polynucleotide, and can precipitate or bind to the corresponding native protein in a cell or tissue preparation or
5 in a cell-free extract of an *in vitro* expression system.

Methods for production of monoclonal and polyclonal antibodies that specifically bind a selected antigen are well known in the art. The antibodies specifically bind to epitopes present in the polypeptides encoded by polynucleotides disclosed in the Sequence Listing. Typically, at least 6, 8, 10, or 12 contiguous amino
10 acids are required to form an epitope. Epitopes that involve non-contiguous amino acids may require a longer polypeptide, *e.g.*, at least 15, 25, or 50 amino acids. Antibodies that specifically bind to human polypeptides encoded by the provided polynucleotides should provide a detection signal at least 5-, 10-, or 20-fold higher than
15 a detection signal provided with other proteins when used in Western blots or other immunochemical assays. Preferably, antibodies that specifically polypeptides of the invention do not bind to other proteins in immunochemical assays at detectable levels and can immunoprecipitate the specific polypeptide from solution.

The invention also contemplates naturally occurring antibodies specific for a polypeptide of the invention. For example, serum antibodies to a polypeptide of
20 the invention in a human population can be purified by methods well known in the art, *e.g.*, by passing antiserum over a column to which the corresponding selected polypeptide or fusion protein is bound. The bound antibodies can then be eluted from the column, for example using a buffer with a high salt concentration.

In addition to the antibodies discussed above, the invention also
25 contemplates genetically engineered antibodies, antibody derivatives (*e.g.*, single chain antibodies, antibody fragments (*e.g.*, Fab, *etc.*)), according to methods well known in the art.

Other embodiments of the present invention include humanized monoclonal antibodies capable of binding to the polypeptides of the invention. The
30 phrase "humanized antibody" refers to an antibody derived from a non-human antibody - typically a mouse monoclonal antibody. Alternatively, a humanized antibody may be derived from a chimeric antibody that retains or substantially retains the antigen-binding properties of the parental, non-human, antibody but which exhibits diminished immunogenicity as compared to the parental antibody when administered to humans.
35 The phrase "chimeric antibody," as used herein, refers to an antibody containing

sequence derived from two different antibodies (*see, e.g.*, U.S. Patent No. 4,816,567) which typically originate from different species. Most typically, chimeric antibodies comprise human and murine antibody fragments, generally human constant and mouse variable regions.

5 Because humanized antibodies are far less immunogenic in humans than the parental mouse monoclonal antibodies, they can be used for the treatment of humans with far less risk of anaphylaxis. Thus, these antibodies may be preferred in therapeutic applications that involve *in vivo* administration to a human such as, *e.g.*, use as radiation sensitizers for the treatment of neoplastic disease or use in methods to reduce the side
10 effects of, *e.g.*, cancer therapy.

 Humanized antibodies may be achieved by a variety of methods including, for example: (1) grafting the non-human complementarity determining regions (CDRs) onto a human framework and constant region (a process referred to in the art as "humanizing"), or, alternatively, (2) transplanting the entire non-human
15 variable domains, but "cloaking" them with a human-like surface by replacement of surface residues (a process referred to in the art as "veneering"). In the present invention, humanized antibodies will include both "humanized" and "veneered" antibodies. These methods are disclosed in, *e.g.*, Jones et al., *Nature* 321:522-525 (1986); Morrison et al., *Proc. Natl. Acad. Sci., U.S.A.*, 81:6851-6855 (1984); Morrison and Oi, *Adv. Immunol.*, 44:65-92 (1988); Verhoever et al., *Science* 239:1534-1536
20 (1988); Padlan, *Molec. Immunol.* 28:489-498 (1991); Padlan, *Molec. Immunol.* 31(3):169-217 (1994); and Kettleborough, C.A. et al., *Protein Eng.* 4(7):773-83 (1991) each of which is incorporated herein by reference.

 The phrase "complementarity determining region" refers to amino acid
25 sequences which together define the binding affinity and specificity of the natural Fv region of a native immunoglobulin binding site. *See, e.g.*, Chothia et al., *J. Mol. Biol.* 196:901-917 (1987); Kabat et al., U.S. Dept. of Health and Human Services NIH Publication No. 91-3242 (1991). The phrase "constant region" refers to the portion of the antibody molecule that confers effector functions. In the present invention, mouse
30 constant regions are substituted by human constant regions. The constant regions of the subject humanized antibodies are derived from human immunoglobulins. The heavy chain constant region can be selected from any of the five isotypes: alpha, delta, epsilon, gamma or mu.

 One method of humanizing antibodies comprises aligning the non-
35 human heavy and light chain sequences to human heavy and light chain sequences,

selecting and replacing the non-human framework with a human framework based on such alignment, molecular modeling to predict the conformation of the humanized sequence and comparing to the conformation of the parent antibody. This process is followed by repeated back mutation of residues in the CDR region which disturb the structure of the CDRs until the predicted conformation of the humanized sequence model closely approximates the conformation of the non-human CDRs of the parent non-human antibody. Such humanized antibodies may be further derivatized to facilitate uptake and clearance, *e.g.*, via Ashwell receptors. *See, e.g.*, U.S. Patent Nos. 5,530,101 and 5,585,089 which patents are incorporated herein by reference.

Humanized antibodies can also be produced using transgenic animals that are engineered to contain human immunoglobulin loci. For example, WO 98/24893 discloses transgenic animals having a human Ig locus wherein the animals do not produce functional endogenous immunoglobulins due to the inactivation of endogenous heavy and light chain loci. WO 91/10741 also discloses transgenic non-primate mammalian hosts capable of mounting an immune response to an immunogen, wherein the antibodies have primate constant and/or variable regions, and wherein the endogenous immunoglobulin-encoding loci are substituted or inactivated. WO 96/30498 discloses the use of the Cre/Lox system to modify the immunoglobulin locus in a mammal, such as to replace all or a portion of the constant or variable region to form a modified antibody molecule. WO 94/02602 discloses non-human mammalian hosts having inactivated endogenous Ig loci and functional human Ig loci. U.S. Patent No. 5,939,598 discloses methods of making transgenic mice in which the mice lack endogenous heavy chains, and express an exogenous immunoglobulin locus comprising one or more xenogeneic constant regions.

Using a transgenic animal described above, an immune response can be produced to a selected antigenic molecule, and antibody-producing cells can be removed from the animal and used to produce hybridomas that secrete human monoclonal antibodies. Immunization protocols, adjuvants, and the like are known in the art, and are used in immunization of, for example, a transgenic mouse as described in WO 96/33735. This publication discloses monoclonal antibodies against a variety of antigenic molecules including IL-6, IL-8, TNF, human CD4, L-selectin, gp39, and tetanus toxin. The monoclonal antibodies can be tested for the ability to inhibit or neutralize the biological activity or physiological effect of the corresponding protein. WO 96/33735 discloses that monoclonal antibodies against IL-8, derived from immune cells of transgenic mice immunized with IL-8, blocked IL-8-induced functions of

neutrophils. Human monoclonal antibodies with specificity for the antigen used to immunize transgenic animals are also disclosed in WO 96/34096.

Polynucleotides or Arrays for Diagnostics

- 5 Polynucleotide arrays are created by spotting polynucleotide probes onto a substrate (e.g., glass, nitrocellulose, *etc.*) in a two-dimensional matrix or array having bound probes. The probes can be bound to the substrate by either covalent bonds or by non-specific interactions, such as hydrophobic interactions. Samples of polynucleotides can be detectably labeled (e.g., using radioactive or fluorescent labels) and then
- 10 hybridized to the probes. Double stranded polynucleotides, comprising the labeled sample polynucleotides bound to probe polynucleotides, can be detected once the unbound portion of the sample is washed away. Techniques for constructing arrays and methods of using these arrays are described in EP 799 897; WO 97/29212; WO 97/27317; EP 785 280; WO 97/02357; U.S. Patent No. 5,593,839; U.S. Patent No.
- 15 5,578,832; EP 728 520; U.S. Patent No. 5,599,695; EP 721 016; U.S. Patent No. 5,556,752; WO 95/22058; and U.S. Patent No. 5,631,734. Arrays can be used to, for example, examine differential expression of genes and can be used to determine gene function. For example, arrays can be used to detect differential expression of a polynucleotide between a test cell and control cell (e.g., cancer cells and normal cells).
- 20 For example, high expression of a particular message in a cancer cell, which is not observed in a corresponding normal cell, can indicate a cancer specific gene product. Exemplary uses of arrays are further described in, for example, Pappalarado et al., *Sem. Radiation Oncol.* (1998) 8:217; and Ramsay, *Nature Biotechnol.* (1998) 16:40.

Differential Expression in Diagnosis

- 25 The polynucleotides of the invention can also be used to detect differences in expression levels between two cells, e.g., as a method to identify abnormal or diseased tissue in a human. For polynucleotides corresponding to profiles of protein families, the choice of tissue can be selected according to the putative biological function. In general, the expression of a gene corresponding to a specific
- 30 polynucleotide is compared between a first tissue that is suspected of being diseased and a second, normal tissue of the human. The tissue suspected of being abnormal or diseased can be derived from a different tissue type of the human, but preferably it is derived from the same tissue type; for example an intestinal polyp or other abnormal growth should be compared with normal intestinal tissue. The normal tissue can be the

same tissue as that of the test sample, or any normal tissue of the patient, especially those that express the polynucleotide-related gene of interest (*e.g.*, brain, thymus, testis, heart, prostate, placenta, spleen, small intestine, skeletal muscle, pancreas, and the mucosal lining of the colon). A difference between the polynucleotide-related gene,
5 mRNA, or protein in the two tissues which are compared, for example in molecular weight, amino acid or nucleotide sequence, or relative abundance, indicates a change in the gene, or a gene which regulates it, in the tissue of the human that was suspected of being diseased. Examples of detection of differential expression and its use in diagnosis of cancer are described in U.S. Patent Nos. 5,688,641 and 5,677,125.

10 A genetic predisposition to disease in a human can also be detected by comparing expression levels of an mRNA or protein corresponding to a polynucleotide of the invention in a fetal tissue with levels associated in normal fetal tissue. Fetal tissues that are used for this purpose include, but are not limited to, amniotic fluid, chorionic villi, blood, and the blastomere of an *in vitro*-fertilized embryo. The
15 comparable normal polynucleotide-related gene is obtained from any tissue. The mRNA or protein is obtained from a normal tissue of a human in which the polynucleotide-related gene is expressed. Differences such as alterations in the nucleotide sequence or size of the same product of the fetal polynucleotide-related gene or mRNA, or alterations in the molecular weight, amino acid sequence, or relative abundance of fetal
20 protein, can indicate a germline mutation in the polynucleotide-related gene of the fetus, which indicates a genetic predisposition to disease. In general, diagnostic, prognostic, and other methods of the invention based on differential expression involve detection of a level or amount of a gene product, particularly a differentially expressed gene product, in a test sample obtained from a patient suspected of having or being susceptible to a
25 disease (*e.g.*, breast cancer, lung cancer, colon cancer and/or metastatic forms thereof), and comparing the detected levels to those levels found in normal cells (*e.g.*, cells substantially unaffected by cancer) and/or other control cells (*e.g.*, to differentiate a cancerous cell from a cell affected by dysplasia). Furthermore, the severity of the disease can be assessed by comparing the detected levels of a differentially expressed
30 gene product with those levels detected in samples representing the levels of differentially gene product associated with varying degrees of severity of disease. It should be noted that use of the term "diagnostic" herein is not necessarily meant to exclude "prognostic" or "prognosis," but rather is used as a matter of convenience.

The term "differentially expressed gene" is generally intended to
35 encompass a polynucleotide that can, for example, include an open reading frame

encoding a gene product (e.g., a polypeptide), and/or introns of such genes and adjacent 5' and 3' non-coding nucleotide sequences involved in the regulation of expression, up to about 20 kb beyond the coding region, but possibly further in either direction. The gene can be introduced into an appropriate vector for extrachromosomal maintenance or for integration into a host genome. In general, a difference in expression level associated with a decrease in expression level of at least about 25%, usually at least about 50% to 75%, more usually at least about 90% or more is indicative of a differentially expressed gene of interest, i.e., a gene that is underexpressed or down-regulated in the test sample relative to a control sample. Furthermore, a difference in expression level associated with an increase in expression of at least about 25%, usually at least about 50% to 75%, more usually at least about 90% and can be at least about 1 1/2-fold, usually at least about 2-fold to about 10-fold, and can be about 100-fold to about 1,000-fold increase relative to a control sample is indicative of a differentially expressed gene of interest, i.e., an overexpressed or up-regulated gene.

"Differentially expressed polynucleotide" as used herein means a nucleic acid molecule (RNA or DNA) comprising a sequence that represents a differentially expressed gene, e.g., the differentially expressed polynucleotide comprises a sequence (e.g., an open reading frame encoding a gene product) that uniquely identifies a differentially expressed gene so that detection of the differentially expressed polynucleotide in a sample is correlated with the presence of a differentially expressed gene in a sample. "Differentially expressed polynucleotides" is also meant to encompass fragments of the disclosed polynucleotides, e.g., fragments retaining biological activity, as well as nucleic acids homologous, substantially similar, or substantially identical (e.g., having about 90% sequence identity) to the disclosed polynucleotides.

"Diagnosis" as used herein generally includes determination of a subject's susceptibility to a disease or disorder, determination as to whether a subject is presently affected by a disease or disorder, as well as to the prognosis of a subject affected by a disease or disorder (e.g., identification of pre-metastatic or metastatic cancerous states, stages of cancer, or responsiveness of cancer to therapy). The present invention particularly encompasses diagnosis of subjects in the context of breast cancer (e.g., carcinoma *in situ* (e.g., ductal carcinoma *in situ*), estrogen receptor (ER)-positive breast cancer, ER-negative breast cancer, or other forms and/or stages of breast cancer), lung cancer (e.g., small cell carcinoma, non-small cell carcinoma, mesothelioma, and

other forms and/or stages of lung cancer), and colon cancer (*e.g.*, adenomatous polyp, colorectal carcinoma, and other forms and/or stages of colon cancer).

“Sample” or “biological sample” as used throughout here are generally meant to refer to samples of biological fluids or tissues, particularly samples obtained from tissues, especially from cells of the type associated with the disease for which the diagnostic application is designed (*e.g.*, ductal adenocarcinoma), and the like. “Samples” is also meant to encompass derivatives and fractions of such samples (*e.g.*, cell lysates). Where the sample is solid tissue, the cells of the tissue can be dissociated or tissue sections can be analyzed.

Methods of the subject invention useful in diagnosis or prognosis typically involve comparison of the abundance of a selected differentially expressed gene product in a sample of interest with that of a control to determine any relative differences in the expression of the gene product, where the difference can be measured qualitatively and/or quantitatively. Quantitation can be accomplished, for example, by comparing the level of expression product detected in the sample with the amounts of product present in a standard curve. A comparison can be made visually; by using a technique such as densitometry, with or without computerized assistance; by preparing a representative library of cDNA clones of mRNA isolated from a test sample, sequencing the clones in the library to determine that number of cDNA clones corresponding to the same gene product, and analyzing the number of clones corresponding to that same gene product relative to the number of clones of the same gene product in a control sample; or by using an array to detect relative levels of hybridization to a selected sequence or set of sequences, and comparing the hybridization pattern to that of a control. The differences in expression are then correlated with the presence or absence of an abnormal expression pattern. A variety of different methods for determining the nucleic acid abundance in a sample are known to those of skill in the art (see, *e.g.*, WO 97/27317). In general, diagnostic assays of the invention involve detection of a gene product of a the polynucleotide sequence (*e.g.*, mRNA or polypeptide) that corresponds to a sequence of SEQ ID NOs:1-3351. The patient from whom the sample is obtained can be apparently healthy, susceptible to disease (*e.g.*, as determined by family history or exposure to certain environmental factors), or can already be identified as having a condition in which altered expression of a gene product of the invention is implicated.

Diagnosis can be determined based on detected gene product expression levels of a gene product encoded by at least one, preferably at least two or more, at least

other forms and/or stages of lung cancer), and colon cancer (*e.g.*, adenomatous polyp, colorectal carcinoma, and other forms and/or stages of colon cancer).

“Sample” or “biological sample” as used throughout here are generally meant to refer to samples of biological fluids or tissues, particularly samples obtained
5 from tissues, especially from cells of the type associated with the disease for which the diagnostic application is designed (*e.g.*, ductal adenocarcinoma), and the like. “Samples” is also meant to encompass derivatives and fractions of such samples (*e.g.*, cell lysates). Where the sample is solid tissue, the cells of the tissue can be dissociated or tissue sections can be analyzed.

10 Methods of the subject invention useful in diagnosis or prognosis typically involve comparison of the abundance of a selected differentially expressed gene product in a sample of interest with that of a control to determine any relative differences in the expression of the gene product, where the difference can be measured qualitatively and/or quantitatively. Quantitation can be accomplished, for example, by
15 comparing the level of expression product detected in the sample with the amounts of product present in a standard curve. A comparison can be made visually; by using a technique such as densitometry, with or without computerized assistance; by preparing a representative library of cDNA clones of mRNA isolated from a test sample, sequencing the clones in the library to determine that number of cDNA clones
20 corresponding to the same gene product, and analyzing the number of clones corresponding to that same gene product relative to the number of clones of the same gene product in a control sample; or by using an array to detect relative levels of hybridization to a selected sequence or set of sequences, and comparing the hybridization pattern to that of a control. The differences in expression are then
25 correlated with the presence or absence of an abnormal expression pattern. A variety of different methods for determining the nucleic acid abundance in a sample are known to those of skill in the art (see, *e.g.*, WO 97/27317). In general, diagnostic assays of the invention involve detection of a gene product of a the polynucleotide sequence (*e.g.*, mRNA or polypeptide) that corresponds to a sequence of SEQ ID NOs:1-3351. The
30 patient from whom the sample is obtained can be apparently healthy, susceptible to disease (*e.g.*, as determined by family history or exposure to certain environmental factors), or can already be identified as having a condition in which altered expression of a gene product of the invention is implicated.

Diagnosis can be determined based on detected gene product expression
35 levels of a gene product encoded by at least one, preferably at least two or more, at least

3 or more, or at least 4 or more of the polynucleotides having a sequence set forth in SEQ ID NOs:1-3351, and can involve detection of expression of genes corresponding to all of SEQ ID NOs:1-3351 and/or additional sequences that can serve as additional diagnostic markers and/or reference sequences. Where the diagnostic method is
5 designed to detect the presence or susceptibility of a patient to cancer, the assay preferably involves detection of a gene product encoded by a gene corresponding to a polynucleotide that is differentially expressed in cancer. Examples of such differentially expressed polynucleotides are described in the Examples below. Given the provided polynucleotides and information regarding their relative expression levels provided
10 herein, assays using such polynucleotides and detection of their expression levels in diagnosis and prognosis will be readily apparent to the ordinarily skilled artisan.

Any of a variety of detectable labels can be used in connection with the various embodiments of the diagnostic methods of the invention. Suitable detectable labels include fluorochromes, (e.g., fluorescein isothiocyanate (FITC), rhodamine,
15 Texas Red, phycoerythrin, allophycocyanin, 6-carboxyfluorescein (6-FAM), 2',7'-dimethoxy-4',5'-dichloro-6-carboxyfluorescein, 6-carboxy-X-rhodamine (ROX), 6-carboxy-2',4',7',4,7-hexachlorofluorescein (HEX), 5-carboxyfluorescein (5-FAM) or N,N,N',N'-tetramethyl-6-carboxyrhodamine (TAMRA)), radioactive labels, (e.g., ^{32}P , ^{35}S , ^3H , etc.), and the like. The detectable label can involve a two stage systems (e.g.,
20 biotin-avidin, hapten-anti-hapten antibody, etc.)

Reagents specific for the polynucleotides and polypeptides of the invention, such as antibodies and nucleotide probes, can be supplied in a kit for detecting the presence of an expression product in a biological sample. The kit can also contain buffers or labeling components, as well as instructions for using the reagents to
25 detect and quantify expression products in the biological sample. Exemplary embodiments of the diagnostic methods of the invention are described below in more detail.

Polypeptide detection in diagnosis. In one embodiment, the test sample is assayed for the level of a differentially expressed polypeptide. Diagnosis can be
30 accomplished using any of a number of methods to determine the absence or presence or altered amounts of the differentially expressed polypeptide in the test sample. For example, detection can utilize staining of cells or histological sections with labeled antibodies, performed in accordance with conventional methods. Cells can be permeabilized to stain cytoplasmic molecules. In general, antibodies that specifically
35 bind a differentially expressed polypeptide of the invention are added to a sample, and

35

incubated for a period of time sufficient to allow binding to the epitope, usually at least about 10 minutes. The antibody can be detectably labeled for direct detection (*e.g.*, using radioisotopes, enzymes, fluorescers, chemilumescers, and the like), or can be used in conjunction with a second stage antibody or reagent to detect binding (*e.g.*,
5 biotin with horseradish peroxidase-conjugated avidin, a secondary antibody conjugated to a fluorescent compound, *e.g.*, fluorescein, rhodamine, Texas red, *etc.*). The absence or presence of antibody binding can be determined by various methods, including flow cytometry of dissociated cells, microscopy, radiography, scintillation counting, *etc.* Any suitable alternative methods can of qualitative or quantitative detection of levels or
10 amounts of differentially expressed polypeptide can be used, for example ELISA, western blot, immunoprecipitation, radioimmunoassay, *etc.*

mRNA detection. The diagnostic methods of the invention can also or alternatively involve detection of mRNA encoded by a gene corresponding to a differentially expressed polynucleotides of the invention. Any suitable qualitative or
15 quantitative methods known in the art for detecting specific mRNAs can be used. mRNA can be detected by, for example, *in situ* hybridization in tissue sections, by reverse transcriptase-PCR, or in Northern blots containing poly A+ mRNA. One of skill in the art can readily use these methods to determine differences in the size or amount of mRNA transcripts between two samples. mRNA expression levels in a
20 sample can also be determined by generation of a library of expressed sequence tags (ESTs) from the sample, where the EST library is representative of sequences present in the sample (Adams, et al., (1991) Science 252:1651). Enumeration of the relative representation of ESTs within the library can be used to approximate the relative representation of the gene transcript within the starting sample. The results of EST
25 analysis of a test sample can then be compared to EST analysis of a reference sample to determine the relative expression levels of a selected polynucleotide, particularly a polynucleotide corresponding to one or more of the differentially expressed genes described herein. Alternatively, gene expression in a test sample can be performed using serial analysis of gene expression (SAGE) methodology (*e.g.*, Velculescu et al.,
30 *Science* (1995) 270:484) or differential display (DD) methodology (see, *e.g.*, U.S. Patent NOs. 5,776,683 and 5,807,680).

Alternatively, gene expression can be analyzed using hybridization analysis. Oligonucleotides or cDNA can be used to selectively identify or capture DNA or RNA of specific sequence composition, and the amount of RNA or cDNA hybridized
35 to a known capture sequence determined qualitatively or quantitatively, to provide

information about the relative representation of a particular message within the pool of cellular messages in a sample. Hybridization analysis can be designed to allow for concurrent screening of the relative expression of hundreds to thousands of genes by using, for example, array-based technologies having high density formats, including
5 filters, microscope slides, or microchips, or solution-based technologies that use spectroscopic analysis (*e.g.*, mass spectrometry). One exemplary use of arrays in the diagnostic methods of the invention is described below in more detail.

Use of a single gene in diagnostic applications. The diagnostic methods of the invention can focus on the expression of a single differentially expressed gene.
10 For example, the diagnostic method can involve detecting a differentially expressed gene, or a polymorphism of such a gene (*e.g.*, a polymorphism in an coding region or control region), that is associated with disease. Disease-associated polymorphisms can include deletion or truncation of the gene, mutations that alter expression level and/or affect activity of the encoded protein, *etc.*

15 A number of methods are available for analyzing nucleic acids for the presence of a specific sequence, *e.g.*, a disease associated polymorphism. Where large amounts of DNA are available, genomic DNA is used directly. Alternatively, the region of interest is cloned into a suitable vector and grown in sufficient quantity for analysis. Cells that express a differentially expressed gene can be used as a source of
20 mRNA, which can be assayed directly or reverse transcribed into cDNA for analysis. The nucleic acid can be amplified by conventional techniques, such as the polymerase chain reaction (PCR), to provide sufficient amounts for analysis, and a detectable label can be included in the amplification reaction (*e.g.*, using a detectably labeled primer or detectably labeled oligonucleotides) to facilitate detection. Alternatively, various
25 methods are also known in the art that utilize oligonucleotide ligation as a means of detecting polymorphisms, see *e.g.*, Riley et al., *Nucl. Acids Res.* (1990) 18:2887; and Delahunty et al., *Am. J. Hum. Genet.* (1996) 58:1239.

The amplified or cloned sample nucleic acid can be analyzed by one of a number of methods known in the art. The nucleic acid can be sequenced by dideoxy or
30 other methods, and the sequence of bases compared to a selected sequence, *e.g.*, to a wild-type sequence. Hybridization with the polymorphic or variant sequence can also be used to determine its presence in a sample (*e.g.*, by Southern blot, dot blot, *etc.*). The hybridization pattern of a polymorphic or variant sequence and a control sequence to an array of oligonucleotide probes immobilized on a solid support, as described in U.S.
35 Patent No. 5,445,934, or in WO 95/35505, can also be used as a means of identifying

polymorphic or variant sequences associated with disease. Single strand conformational polymorphism (SSCP) analysis, denaturing gradient gel electrophoresis (DGGE), and heteroduplex analysis in gel matrices are used to detect conformational changes created by DNA sequence variation as alterations in electrophoretic mobility.

5 Alternatively, where a polymorphism creates or destroys a recognition site for a restriction endonuclease, the sample is digested with that endonuclease, and the products size fractionated to determine whether the fragment was digested. Fractionation is performed by gel or capillary electrophoresis, particularly acrylamide or agarose gels.

10 Screening for mutations in a gene can be based on the functional or antigenic characteristics of the protein. Protein truncation assays are useful in detecting deletions that can affect the biological activity of the protein. Various immunoassays designed to detect polymorphisms in proteins can be used in screening. Where many diverse genetic mutations lead to a particular disease phenotype, functional protein
15 assays have proven to be effective screening tools. The activity of the encoded protein can be determined by comparison with the wild-type protein.

Pattern matching in diagnosis using arrays. In another embodiment, the diagnostic and/or prognostic methods of the invention involve detection of expression of a selected set of genes in a test sample to produce a test expression pattern (TEP).
20 The TEP is compared to a reference expression pattern (REP), which is generated by detection of expression of the selected set of genes in a reference sample (*e.g.*, a positive or negative control sample). The selected set of genes includes at least one of the genes of the invention, which genes correspond to the polynucleotide sequences of SEQ ID NOs:1-3351. Of particular interest is a selected set of genes that includes genes
25 differentially expressed in the disease for which the test sample is to be screened.

"Reference sequences" or "reference polynucleotides" as used herein in the context of differential gene expression analysis and diagnosis/prognosis refers to a selected set of polynucleotides, which selected set includes at least one or more of the differentially expressed polynucleotides described herein. A plurality of reference
30 sequences, preferably comprising positive and negative control sequences, can be included as reference sequences. Additional suitable reference sequences are found in Genbank, Unigene, and other nucleotide sequence databases (including, *e.g.*, expressed sequence tag (EST), partial, and full-length sequences).

"Reference array" means an array having reference sequences for use in
35 hybridization with a sample, where the reference sequences include all, at least one of,

or any subset of the differentially expressed polynucleotides described herein. Usually such an array will include at least 3 different reference sequences, and can include any one or all of the provided differentially expressed sequences. Arrays of interest can further comprise sequences, including polymorphisms, of other genetic sequences, particularly other sequences of interest for screening for a disease or disorder (e.g., cancer, dysplasia, or other related or unrelated diseases, disorders, or conditions). The oligonucleotide sequence on the array will usually be at least about 12 nt in length, and can be of about the length of the provided sequences, or can extend into the flanking regions to generate fragments of 100 nt to 200 nt in length or more. Reference arrays can be produced according to any suitable methods known in the art. For example, methods of producing large arrays of oligonucleotides are described in U.S. Patent NOs. 5,134,854 and 5,445,934 using light-directed synthesis techniques. Using a computer controlled system, a heterogeneous array of monomers is converted, through simultaneous coupling at a number of reaction sites, into a heterogeneous array of polymers. Alternatively, microarrays are generated by deposition of pre-synthesized oligonucleotides onto a solid substrate, for example as described in PCT published application no. WO 95/35505.

A "reference expression pattern" or "REP" as used herein refers to the relative levels of expression of a selected set of genes, particularly of differentially expressed genes, that is associated with a selected cell type, e.g., a normal cell, a cancerous cell, a cell exposed to an environmental stimulus, and the like. A "test expression pattern" or "TEP" refers to relative levels of expression of a selected set of genes, particularly of differentially expressed genes, in a test sample (e.g., a cell of unknown or suspected disease state, from which mRNA is isolated).

REPs can be generated in a variety of ways according to methods well known in the art. For example, REPs can be generated by hybridizing a control sample to an array having a selected set of polynucleotides (particularly a selected set of differentially expressed polynucleotides), acquiring the hybridization data from the array, and storing the data in a format that allows for ready comparison of the REP with a TEP. Alternatively, all expressed sequences in a control sample can be isolated and sequenced, e.g., by isolating mRNA from a control sample, converting the mRNA into cDNA, and sequencing the cDNA. The resulting sequence information roughly or precisely reflects the identity and relative number of expressed sequences in the sample. The sequence information can then be stored in a format (e.g., a computer-readable format) that allows for ready comparison of the REP with a TEP. The REP can be

normalized prior to or after data storage, and/or can be processed to selectively remove sequences of expressed genes that are of less interest or that might complicate analysis (e.g., some or all of the sequences associated with housekeeping genes can be eliminated from REP data).

- 5 TEPs can be generated in a manner similar to REPs, e.g., by hybridizing a test sample to an array having a selected set of polynucleotides, particularly a selected set of differentially expressed polynucleotides, acquiring the hybridization data from the array, and storing the data in a format that allows for ready comparison of the TEP with a REP. The REP and TEP to be used in a comparison can be generated simultaneously,
10 or the TEP can be compared to previously generated and stored REPs.

- In one embodiment of the invention, comparison of a TEP with a REP involves hybridizing a test sample with a reference array, where the reference array has one or more reference sequences for use in hybridization with a sample. The reference sequences include all, at least one of, or any subset of the differentially expressed
15 polynucleotides described herein. Hybridization data for the test sample is acquired, the data normalized, and the produced TEP compared with a REP generated using an array having the same or similar selected set of differentially expressed polynucleotides. Probes that correspond to sequences differentially expressed between the two samples will show decreased or increased hybridization efficiency for one of the samples
20 relative to the other.

- Methods for collection of data from hybridization of samples with a reference arrays are well known in the art. For example, the polynucleotides of the reference and test samples can be generated using a detectable fluorescent label, and hybridization of the polynucleotides in the samples detected by scanning the
25 microarrays for the presence of the detectable label using, for example, a microscope and light source for directing light at a substrate. A photon counter detects fluorescence from the substrate, while an x-y translation stage varies the location of the substrate. A confocal detection device that can be used in the subject methods is described in U.S. Patent No. 5,631,734. A scanning laser microscope is described in Shalon et al.,
30 *Genome Res.* (1996) 6:639. A scan, using the appropriate excitation line, is performed for each fluorophore used. The digital images generated from the scan are then combined for subsequent analysis. For any particular array element, the ratio of the fluorescent signal from one sample (e.g., a test sample) is compared to the fluorescent signal from another sample (e.g., a reference sample), and the relative signal intensity
35 determined.

Methods for analyzing the data collected from hybridization to arrays are well known in the art. For example, where detection of hybridization involves a fluorescent label, data analysis can include the steps of determining fluorescent intensity as a function of substrate position from the data collected, removing outliers, *i.e.*, data
5 deviating from a predetermined statistical distribution, and calculating the relative binding affinity of the targets from the remaining data. The resulting data can be displayed as an image with the intensity in each region varying according to the binding affinity between targets and probes.

In general, the test sample is classified as having a gene expression
10 profile corresponding to that associated with a disease or non-disease state by comparing the TEP generated from the test sample to one or more REPs generated from reference samples (*e.g.*, from samples associated with cancer or specific stages of cancer, dysplasia, samples affected by a disease other than cancer, normal samples, *etc.*). The criteria for a match or a substantial match between a TEP and a REP include
15 expression of the same or substantially the same set of reference genes, as well as expression of these reference genes at substantially the same levels (*e.g.*, no significant difference between the samples for a signal associated with a selected reference sequence after normalization of the samples, or at least no greater than about 25% to about 40% difference in signal strength for a given reference sequence. In general, a
20 pattern match between a TEP and a REP includes a match in expression, preferably a match in qualitative or quantitative expression level, of at least one of, all or any subset of the differentially expressed genes of the invention.

Pattern matching can be performed manually, or can be performed using a computer program. Methods for preparation of substrate matrices (*e.g.*, arrays),
25 design of oligonucleotides for use with such matrices, labeling of probes, hybridization conditions, scanning of hybridized matrices, and analysis of patterns generated, including comparison analysis, are described in, for example, U.S. Patent No. 5,800,992.

Diagnosis, Prognosis and Management of Cancer

30 The polynucleotides of the invention and their gene products are of particular interest as genetic or biochemical markers (*e.g.*, in blood or tissues) that will detect the earliest changes along the carcinogenesis pathway and/or to monitor the efficacy of various therapies and preventive interventions. For example, the level of expression of certain polynucleotides can be indicative of a poorer prognosis, and

therefore warrant more aggressive chemo- or radio-therapy for a patient or vice versa. The correlation of novel surrogate tumor specific features with response to treatment and outcome in patients can define prognostic indicators that allow the design of tailored therapy based on the molecular profile of the tumor. These therapies include antibody targeting and gene therapy. Determining expression of certain polynucleotides and comparison of a patient's profile with known expression in normal tissue and variants of the disease allows a determination of the best possible treatment for a patient, both in terms of specificity of treatment and in terms of comfort level of the patient. Surrogate tumor markers, such as polynucleotide expression, can also be used to better classify, and thus diagnose and treat, different forms and disease states of cancer. Two classifications widely used in oncology that can benefit from identification of the expression levels of the polynucleotides of the invention are staging of the cancerous disorder, and grading the nature of the cancerous tissue.

The polynucleotides of the invention can be useful to monitor patients having or susceptible to cancer to detect potentially malignant events at a molecular level before they are detectable at a gross morphological level. Furthermore, a polynucleotide of the invention identified as important for one type of cancer can also have implications for development or risk of development of other types of cancer, *e.g.*, where a polynucleotide is differentially expressed across various cancer types. Thus, for example, expression of a polynucleotide that has clinical implications for metastatic colon cancer can also have clinical implications for stomach cancer or endometrial cancer.

Staging. Staging is a process used by physicians to describe how advanced the cancerous state is in a patient. Generally, if a cancer is only detectable in the area of the primary lesion without having spread to any lymph nodes it is called Stage I. If it has spread only to the closest lymph nodes, it is called Stage II. In Stage III, the cancer has generally spread to the lymph nodes in near proximity to the site of the primary lesion. Cancers that have spread to a distant part of the body, such as the liver, bone, brain or other site, are Stage IV, the most advanced stage.

The polynucleotides of the invention can facilitate fine-tuning of the staging process by identifying markers for the aggressivity of a cancer, *e.g.*, the metastatic potential, as well as the presence in different areas of the body. Thus, a Stage II cancer with a polynucleotide signifying a high metastatic potential cancer can be used to change a borderline Stage II tumor to a Stage III tumor, justifying more aggressive

therapy. Conversely, the presence of a polynucleotide signifying a lower metastatic potential allows more conservative staging of a tumor.

Grading of cancers. Grade is a term used to describe how closely a tumor resembles normal tissue of its same type. The microscopic appearance of a tumor is used to identify tumor grade based on parameters such as cell morphology, cellular organization, and other markers of differentiation. As a general rule, the grade of a tumor corresponds to its rate of growth or aggressiveness, with undifferentiated or high-grade tumors being more aggressive than well differentiated or low-grade tumors. The following guidelines are generally used for grading tumors: 1) GX Grade cannot be assessed; 2) G1 Well differentiated; G2 Moderately well differentiated; 3) G3 Poorly differentiated; 4) G4 Undifferentiated. The polynucleotides of the invention can be especially valuable in determining the grade of the tumor, as they not only can aid in determining the differentiation status of the cells of a tumor, they can also identify factors other than differentiation that are valuable in determining the aggressivity of a tumor, such as metastatic potential.

Detection of lung cancer. The polynucleotides of the invention can be used to detect lung cancer in a subject. Although there are more than a dozen different kinds of lung cancer, the two main types of lung cancer are small cell and nonsmall cell, which encompass about 90% of all lung cancer cases. Small cell carcinoma (also called oat cell carcinoma) usually starts in one of the larger bronchial tubes, grows fairly rapidly, and is likely to be large by the time of diagnosis. Nonsmall cell lung cancer (NSCLC) is made up of three general subtypes of lung cancer. Epidermoid carcinoma (also called squamous cell carcinoma) usually starts in one of the larger bronchial tubes and grows relatively slowly. The size of these tumors can range from very small to quite large. Adenocarcinoma starts growing near the outside surface of the lung and can vary in both size and growth rate. Some slowly growing adenocarcinomas are described as alveolar cell cancer. Large cell carcinoma starts near the surface of the lung, grows rapidly, and the growth is usually fairly large when diagnosed. Other less common forms of lung cancer are carcinoid, cylindroma, mucoepidermoid, and malignant mesothelioma.

The polynucleotides of the invention, *e.g.*, polynucleotides differentially expressed in normal cells versus cancerous lung cells (*e.g.*, tumor cells of high or low metastatic potential) or between types of cancerous lung cells (*e.g.*, high metastatic versus low metastatic), can be used to distinguish types of lung cancer as well as identifying traits specific to a certain patient's cancer and selecting an appropriate

therapy. For example, if the patient's biopsy expresses a polynucleotide that is associated with a low metastatic potential, it may justify leaving a larger portion of the patient's lung in surgery to remove the lesion. Alternatively, a smaller lesion with expression of a polynucleotide that is associated with high metastatic potential may justify a more radical removal of lung tissue and/or the surrounding lymph nodes, even if no metastasis can be identified through pathological examination.

Detection of breast cancer. The majority of breast cancers are adenocarcinomas subtypes, which can be summarized as follows: 1) ductal carcinoma *in situ* (DCIS), including comedocarcinoma; 2) infiltrating (or invasive) ductal carcinoma (IDC); 3) lobular carcinoma *in situ* (LCIS); 4) infiltrating (or invasive) lobular carcinoma (ILC); 5) inflammatory breast cancer; 6) medullary carcinoma; 7) mucinous carcinoma; 8) Paget's disease of the nipple; 9) Phyllodes tumor; and 10) tubular carcinoma.

The expression of polynucleotides of the invention can be used in the diagnosis and management of breast cancer, as well as to distinguish between types of breast cancer. Detection of breast cancer can be determined using expression levels of any of the appropriate polynucleotides of the invention, either alone or in combination. Determination of the aggressive nature and/or the metastatic potential of a breast cancer can also be determined by comparing levels of one or more polynucleotides of the invention and comparing levels of another sequence known to vary in cancerous tissue, *e.g.*, ER expression. In addition, development of breast cancer can be detected by examining the ratio of expression of a differentially expressed polynucleotide to the levels of steroid hormones (*e.g.*, testosterone or estrogen) or to other hormones (*e.g.*, growth hormone, insulin). Thus expression of specific marker polynucleotides can be used to discriminate between normal and cancerous breast tissue, to discriminate between breast cancers with different cells of origin, to discriminate between breast cancers with different potential metastatic rates, *etc.*

Detection of colon cancer. The polynucleotides of the invention exhibiting the appropriate expression pattern can be used to detect colon cancer in a subject. Colorectal cancer is one of the most common neoplasms in humans and perhaps the most frequent form of hereditary neoplasia. Prevention and early detection are key factors in controlling and curing colorectal cancer. Colorectal cancer begins as polyps, which are small, benign growths of cells that form on the inner lining of the colon. Over a period of several years, some of these polyps accumulate additional mutations and become cancerous. Multiple familial colorectal cancer disorders have

been identified, which are summarized as follows: 1) Familial adenomatous polyposis (FAP); 2) Gardner's syndrome; 3) Hereditary nonpolyposis colon cancer (HNPCC); and 4) Familial colorectal cancer in Ashkenazi Jews. The expression of appropriate polynucleotides of the invention can be used in the diagnosis, prognosis and management of colorectal cancer. Detection of colon cancer can be determined using expression levels of any of these sequences alone or in combination with the levels of expression. Determination of the aggressive nature and/or the metastatic potential of a colon cancer can be determined by comparing levels of one or more polynucleotides of the invention and comparing total levels of another sequence known to vary in cancerous tissue, *e.g.*, expression of p53, DCC ras, or FAP (see, *e.g.*, Fearon ER, et al., *Cell* (1990) 61(5):759; Hamilton SR et al., *Cancer* (1993) 72:957; Bodmer W, et al., *Nat Genet.* (1994) 4(3):217; Fearon ER, *Ann N Y Acad Sci.* (1995) 768:101). For example, development of colon cancer can be detected by examining the ratio of any of the polynucleotides of the invention to the levels of oncogenes (*e.g.*, ras) or tumor suppressor genes (*e.g.*, FAP or p53). Thus expression of specific marker polynucleotides can be used to discriminate between normal and cancerous colon tissue, to discriminate between colon cancers with different cells of origin, to discriminate between colon cancers with different potential metastatic rates, *etc.*

Use of Polynucleotides to Screen for Peptide Analogs and Antagonists

Polypeptides encoded by the instant polynucleotides and corresponding full length genes can be used to screen peptide libraries to identify binding partners, such as receptors, from among the encoded polypeptides. Peptide libraries can be synthesized according to methods known in the art (see, *e.g.*, U.S. Patent No. 5,010,175, and WO 91/17823). Agonists or antagonists of the polypeptides of the invention can be screened using any available method known in the art, such as signal transduction, antibody binding, receptor binding, mitogenic assays, chemotaxis assays, *etc.* The assay conditions ideally should resemble the conditions under which the native activity is exhibited *in vivo*, that is, under physiologic pH, temperature, and ionic strength. Suitable agonists or antagonists will exhibit strong inhibition or enhancement of the native activity at concentrations that do not cause toxic side effects in the subject. Agonists or antagonists that compete for binding to the native polypeptide can require concentrations equal to or greater than the native concentration, while inhibitors capable of binding irreversibly to the polypeptide can be added in concentrations on the order of the native concentration.

Such screening and experimentation can lead to identification of a novel polypeptide binding partner, such as a receptor, encoded by a gene or a cDNA corresponding to a polynucleotide of the invention, and at least one peptide agonist or antagonist of the novel binding partner. Such agonists and antagonists can be used to modulate, enhance, or inhibit receptor function in cells to which the receptor is native, or in cells that possess the receptor as a result of genetic engineering. Further, if the novel receptor shares biologically important characteristics with a known receptor, information about agonist/antagonist binding can facilitate development of improved agonists/antagonists of the known receptor.

10 Pharmaceutical Compositions and Therapeutic Uses

Pharmaceutical compositions of the invention can comprise polypeptides, antibodies, or polynucleotides (including antisense nucleotides and ribozymes) of the claimed invention in a therapeutically effective amount. The term "therapeutically effective amount" as used herein refers to an amount of a therapeutic agent to treat, ameliorate, or prevent a desired disease or condition, or to exhibit a detectable therapeutic or preventative effect. The effect can be detected by, for example, chemical markers or antigen levels. Therapeutic effects also include reduction in physical symptoms, such as decreased body temperature. The precise effective amount for a subject will depend upon the subject's size and health, the nature and extent of the condition, and the therapeutics or combination of therapeutics selected for administration. Thus, it is not useful to specify an exact effective amount in advance. However, the effective amount for a given situation is determined by routine experimentation and is within the judgment of the clinician. For purposes of the present invention, an effective dose will generally be from about 0.01 mg/kg to 50 mg/kg or 0.05 mg/kg to about 10 mg/kg of the DNA constructs in the individual to which it is administered.

A pharmaceutical composition can also contain a pharmaceutically acceptable carrier. The term "pharmaceutically acceptable carrier" refers to a carrier for administration of a therapeutic agent, such as antibodies or a polypeptide, genes, and other therapeutic agents. The term refers to any pharmaceutical carrier that does not itself induce the production of antibodies harmful to the individual receiving the composition, and which can be administered without undue toxicity. Suitable carriers can be large, slowly metabolized macromolecules such as proteins, polysaccharides, polylactic acids, polyglycolic acids, polymeric amino acids, amino acid copolymers,

and inactive virus particles. Such carriers are well known to those of ordinary skill in the art. Pharmaceutically acceptable carriers in therapeutic compositions can include liquids such as water, saline, glycerol and ethanol. Auxiliary substances, such as wetting or emulsifying agents, pH buffering substances, and the like, can also be present in such vehicles. Typically, the therapeutic compositions are prepared as injectables, either as liquid solutions or suspensions; solid forms suitable for solution in, or suspension in, liquid vehicles prior to injection can also be prepared. Liposomes are included within the definition of a pharmaceutically acceptable carrier. Pharmaceutically acceptable salts can also be present in the pharmaceutical composition, e.g., mineral acid salts such as hydrochlorides, hydrobromides, phosphates, sulfates, and the like; and the salts of organic acids such as acetates, propionates, malonates, benzoates, and the like. A thorough discussion of pharmaceutically acceptable excipients is available in *Remington's Pharmaceutical Sciences* (Mack Pub. Co., New Jersey, 1991).

Delivery Methods. Once formulated, the compositions of the invention can be (1) administered directly to the subject (e.g., as polynucleotide or polypeptides); or (2) delivered *ex vivo*, to cells derived from the subject (e.g., as in *ex vivo* gene therapy). Direct delivery of the compositions will generally be accomplished by parenteral injection, e.g., subcutaneously, intraperitoneally, intravenously or intramuscularly, intratumoral or to the interstitial space of a tissue. Other modes of administration include oral and pulmonary administration, suppositories, and transdermal applications, needles, and gene guns or hyposprays. Dosage treatment can be a single dose schedule or a multiple dose schedule.

Methods for the *ex vivo* delivery and reimplantation of transformed cells into a subject are known in the art and described in e.g., International Publication No. WO 93/14778. Examples of cells useful in *ex vivo* applications include, for example, stem cells, particularly hematopoietic, lymph cells, macrophages, dendritic cells, or tumor cells. Generally, delivery of nucleic acids for both *ex vivo* and *in vitro* applications can be accomplished by, for example, dextran-mediated transfection, calcium phosphate precipitation, polybrene mediated transfection, protoplast fusion, electroporation, encapsulation of the polynucleotide(s) in liposomes, and direct microinjection of the DNA into nuclei, all well known in the art.

Once a gene corresponding to a polynucleotide of the invention has been found to correlate with a proliferative disorder, such as neoplasia, dysplasia, and hyperplasia, the disorder can be amenable to treatment by administration of a

therapeutic agent based on the provided polynucleotide, corresponding polypeptide or other corresponding molecule (e.g., antisense, ribozyme, etc.).

The dose and the means of administration of the inventive pharmaceutical compositions are determined based on the specific qualities of the therapeutic composition, the condition, age, and weight of the patient, the progression of the disease, and other relevant factors. For example, administration of polynucleotide therapeutic compositions agents of the invention includes local or systemic administration, including injection, oral administration, particle gun or catheterized administration, and topical administration. Preferably, the therapeutic polynucleotide composition contains an expression construct comprising a promoter operably linked to a polynucleotide of at least 12, 22, 25, 30, or 35 contiguous nt of the polynucleotide disclosed herein. Various methods can be used to administer the therapeutic composition directly to a specific site in the body. For example, a small metastatic lesion is located and the therapeutic composition injected several times in several different locations within the body of tumor. Alternatively, arteries which serve a tumor are identified, and the therapeutic composition injected into such an artery, in order to deliver the composition directly into the tumor. A tumor that has a necrotic center is aspirated and the composition injected directly into the now empty center of the tumor. The antisense composition is directly administered to the surface of the tumor, for example, by topical application of the composition. X-ray imaging is used to assist in certain of the above delivery methods.

Receptor-mediated targeted delivery of therapeutic compositions containing an antisense polynucleotide, subgenomic polynucleotides, or antibodies to specific tissues can also be used. Receptor-mediated DNA delivery techniques are described in, for example, Findeis et al., *Trends Biotechnol.* (1993) 11:202; Chiou et al., *Gene Therapeutics: Methods And Applications Of Direct Gene Transfer* (J.A. Wolff, ed.) (1994); Wu et al., *J. Biol. Chem.* (1988) 263:621; Wu et al., *J. Biol. Chem.* (1994) 269:542; Zenke et al., *Proc. Natl. Acad. Sci. (USA)* (1990) 87:3655; Wu et al., *J. Biol. Chem.* (1991) 266:338. Therapeutic compositions containing a polynucleotide are administered in a range of about 100 ng to about 200 mg of DNA for local administration in a gene therapy protocol. Concentration ranges of about 500 ng to about 50 mg, about 1 mg to about 2 mg, about 5 mg to about 500 mg, and about 20 mg to about 100 mg of DNA can also be used during a gene therapy protocol. Factors such as method of action (e.g., for enhancing or inhibiting levels of the encoded gene product) and efficacy of transformation and expression are considerations which will

affect the dosage required for ultimate efficacy of the antisense subgenomic polynucleotides. Where greater expression is desired over a larger area of tissue, larger amounts of antisense subgenomic polynucleotides or the same amounts readministered in a successive protocol of administrations, or several administrations to different adjacent or close tissue portions of, for example, a tumor site, may be required to effect a positive therapeutic outcome. In all cases, routine experimentation in clinical trials will determine specific ranges for optimal therapeutic effect. For polynucleotide-related genes encoding polypeptides or proteins with anti-inflammatory activity, suitable use, doses, and administration are described in U.S. Patent No. 5,654,173.

10 The therapeutic polynucleotides and polypeptides of the present invention can be delivered using gene delivery vehicles. The gene delivery vehicle can be of viral or non-viral origin (see generally, Jolly, *Cancer Gene Therapy* (1994) 1:51; Kimura, *Human Gene Therapy* (1994) 5:845; Connelly, *Human Gene Therapy* (1995) 1:185; and Kaplitt, *Nature Genetics* (1994) 6:148). Expression of such coding sequences can be induced using endogenous mammalian or heterologous promoters. Expression of the coding sequence can be either constitutive or regulated.

Viral-based vectors for delivery of a desired polynucleotide and expression in a desired cell are well known in the art. Exemplary viral-based vehicles include, but are not limited to, recombinant retroviruses (see, e.g., WO 90/07936; WO 20 94/03622; WO 93/25698; WO 93/25234; U.S. Patent No. 5, 219,740; WO 93/11230; WO 93/10218; U.S. Patent No. 4,777,127; GB Patent No. 2,200,651; EP 0 345 242; and WO 91/02805), alphavirus-based vectors (e.g., Sindbis virus vectors, Semliki forest virus (ATCC VR-67; ATCC VR-1247), Ross River virus (ATCC VR-373; ATCC VR-1246) and Venezuelan equine encephalitis virus (ATCC VR-923; ATCC VR-1250; 25 ATCC VR 1249; ATCC VR-532), and adeno-associated virus (AAV) vectors (see, e.g., WO 94/12649, WO 93/03769; WO 93/19191; WO 94/28938; WO 95/11984 and WO 95/00655). Administration of DNA linked to killed adenovirus as described in Curiel, *Hum. Gene Ther.* (1992) 3:147 can also be employed.

Non-viral delivery vehicles and methods can also be employed, including, but not limited to, polycationic condensed DNA linked or unlinked to killed adenovirus alone (see, e.g., Curiel, *Hum. Gene Ther.* (1992) 3:147); ligand-linked DNA (see, e.g., Wu, *J. Biol. Chem.* 264:16985 (1989)); eukaryotic cell delivery vehicles cells (see, e.g., U.S. Patent No. 5,814,482; WO 95/07994; WO 96/17072; WO 95/30763; and WO 97/42338) and nucleic charge neutralization or fusion with cell 35 membranes. Naked DNA can also be employed. Exemplary naked DNA introduction

methods are described in WO 90/11092 and U.S. Patent No. 5,580,859. Liposomes that can act as gene delivery vehicles are described in U.S. Patent No. 5,422,120; WO 95/13796; WO 94/23697; WO 91/14445; and EP 0524968. Additional approaches are described in Philip, *Mol. Cell Biol.* 14:2411 (1994), and in Woffendin, *Proc. Natl. Acad. Sci.* (1994) 91:1581.

Further non-viral delivery suitable for use includes mechanical delivery systems such as the approach described in Woffendin et al., *Proc. Natl. Acad. Sci. USA* 91(24):11581 (1994). Moreover, the coding sequence and the product of expression of such can be delivered through deposition of photopolymerized hydrogel materials or use of ionizing radiation (see, e.g., U.S. Patent No. 5,206,152 and WO 92/11033). Other conventional methods for gene delivery that can be used for delivery of the coding sequence include, for example, use of hand-held gene transfer particle gun (see, e.g., U.S. Patent No. 5,149,655); use of ionizing radiation for activating transferred gene (see, e.g., U.S. Patent No. 5,206,152 and WO 92/11033).

The present invention will now be illustrated by reference to the following examples which set forth particularly advantageous embodiments. However, it should be noted that these embodiments are illustrative and are not to be construed as restricting the invention in any way.

EXAMPLES

EXAMPLE 1

SOURCE OF BIOLOGICAL MATERIALS AND OVERVIEW OF NOVEL POLYNUCLEOTIDES EXPRESSED BY THE BIOLOGICAL MATERIALS

5

Cell lines and human normal and tumor tissue were used to construct cDNA libraries from mRNA isolated from the cells and tissues. Most sequences were about 275-300 nucleotides in length. The cells lines include Km12L4-A cell line, a high metastatic colon cancer cell line (Morika, W. A. K. et al., *Cancer Research* (1988) 48:6863). The KM12L4-A cell line is derived from the KM12C cell line. The KM12C cell line, which is poorly metastatic (low metastatic) was established in culture from a Dukes' stage B2 surgical specimen (Morikawa et al. *Cancer Res.* (1988) 48:6863). The KML4-A is a highly metastatic subline derived from KM12C (Yeatman et al. *Nucl. Acids. Res.* (1995) 23:4007; Bao-Ling et al. *Proc. Annu. Meet. Am. Assoc. Cancer. Res.* (1995) 21:3269). The KM12C and KM12C-derived cell lines (e.g., KM12L4, KM12L4-A, etc.) are well-recognized in the art as model cell lines for the study of colon cancer (see, e.g., Moriakawa et al., *supra*; Radinsky et al. *Clin. Cancer Res.* (1995) 1:19; Yeatman et al., (1995) *supra*; Yeatman et al., *Clin. Exp. Metastasis* (1996) 14:246). These and other cell lines and tissue are described in Table 6.

20

The sequences of the isolated polynucleotides were first masked to eliminate low complexity sequences using the XBLAST masking program (Claverie "Effective Large-Scale Sequence Similarity Searches," In: Computer Methods for Macromolecular Sequence Analysis, Doolittle, ed., *Meth. Enzymol.* 266:212-227 Academic Press, NY, NY (1996); see particularly Claverie, in "Automated DNA Sequencing and Analysis Techniques" Adams et al., eds., Chap. 36, p. 267 Academic Press, San Diego, 1994 and Claverie et al. *Comput. Chem.* (1993) 17:191). Generally, masking does not influence the final search results, except to eliminate sequences of relative little interest due to their low complexity, and to eliminate multiple "hits" based on similarity to repetitive regions common to multiple sequences, e.g., Alu repeats. The sequences remaining after masking were then used in a BLASTN vs. Genbank search; sequences that exhibited greater than 70% overlap, 99% identity, and a p value of less than 1×10^{-40} were discarded. Sequences from this search also were discarded if the inclusive parameters were met, but the sequence was ribosomal or vector-derived.

25

30

The resulting sequences from the previous search were classified into three groups (1, 2 and 3 below) and searched in a BLASTX vs. NRP (non-redundant proteins) database search: (1) unknown (no hits in the Genbank search), (2) weak similarity (greater than 45% identity and p value of less than 1×10^{-5}), and (3) high similarity (greater than 60% overlap, greater than 80% identity, and p value less than 1×10^{-5}). Sequences having greater than 70% overlap, greater than 99% identity, and p value of less than 1×10^{-40} were discarded.

The remaining sequences were classified as unknown (no hits), weak similarity, and high similarity (parameters as above). Two searches were performed on these sequences. First, a BLAST vs. EST database search was performed and sequences with greater than 99% overlap, greater than 99% similarity and a p value of less than 1×10^{-40} were discarded. Sequences with a p value of less than 1×10^{-65} when compared to a database sequence of human origin were also excluded. Second, a BLASTN vs. Patent GeneSeq database was performed and sequences having greater than 99% identity, p value less than 1×10^{-40} , and greater than 99% overlap were discarded.

The remaining sequences were subjected to screening using other rules and redundancies in the dataset. Sequences with a p value of less than 1×10^{-111} in relation to a database sequence of human origin were specifically excluded. The final result provided the 3351 sequences listed in the accompanying Sequence Listing. Each identified polynucleotide represents sequence from at least a partial mRNA transcript. Polynucleotides that were determined to be novel were assigned a sequence identification number.

The novel polynucleotides were assigned sequence identification numbers SEQ ID NOs:1-3351. The first 1847 DNA sequences corresponding to the novel polynucleotides are provided in the Sequence Listing in Table 1. DNA sequences corresponding to the novel polynucleotides of SEQ ID NOs:1848-3351 are provided in the Sequence Listing in Table 2. The DNA sequences of Table 2, while numbered SEQ ID 1-1504, correspond to SEQ ID NOs:1848-3351 in the Sequence Listing, *e.g.*, Table 2 SEQ ID 1 is SEQ ID NO:1848, Table 2 SEQ ID 2 is SEQ ID NO:1849; *etc.* Each DNA sequence in Table 4 is uniquely identified by a number that is 1847 less than its SEQ ID NO in the Sequence Listing. Tables 1 and 2 provide: 1) the SEQ ID NO assigned to each sequence for use in the present specification or a corresponding number; 2) the sequence name used as an internal identifier of the sequence; 3) the name assigned to the clone from which the

sequence was isolated; and 4) the number of the cluster to which the sequence is assigned (Cluster ID; where the cluster ID is 0, the sequence was not assigned to any cluster).

Because the provided polynucleotides represent partial mRNA transcripts, two or more polynucleotides of the invention may represent different
5 regions of the same mRNA transcript and the same gene. Thus, if two or more SEQ ID NOs: are identified as belonging to the same clone, then either sequence can be used to obtain the full-length mRNA or gene.

EXAMPLE 2

RESULTS OF PUBLIC DATABASE SEARCH TO IDENTIFY FUNCTION OF GENE PRODUCTS

10

SEQ ID NOs:1-3351 were translated in all three reading frames to determine the best alignment with the individual sequences. These amino acid sequences and nucleotide sequences are referred to, generally, as query sequences, which are aligned with the individual sequences. Query and individual sequences were
15 aligned using the BLAST programs, available over the world wide web at <http://www.ncbi.nlm.nih.gov/BLAST/>. Again the sequences were masked to various extents to prevent searching of repetitive sequences or poly-A sequences, using the XBLAST program for masking low complexity as described above in Example 1.

Tables 3 and 4 (inserted before the claims) show the results of the
20 alignments. Table 3 contains alignment information for SEQ ID NOs:1-1847 and Table 4 contains alignment information for SEQ ID NOs:1848-3351. The DNA sequences of Table 4, while numbered SEQ ID 1-1504, correspond to SEQ ID NOs:1848-3351. Each DNA sequence in Table 4 is uniquely identified by a number that is 1847 less than its SEQ ID NO. Tables 3 and 4 refer to each sequence by its SEQ ID NO or a corresponding number,
25 the accession numbers and descriptions of nearest neighbors from the Genbank and Non-Redundant Protein searches, and the p values of the search results.

For each of SEQ ID NOs:1-1847, the best alignment to a protein or DNA sequence is included in Table 3, and the best alignment for each of SEQ ID NOs:1848-3351 is included in Table 4. The activity of the polypeptide encoded by SEQ ID
30 NOs:1-3351 is the same or similar to the nearest neighbor reported in Table 3 or 4. The accession number of the nearest neighbor is reported, providing a reference to the activities exhibited by the nearest neighbor. The search program and database used for the alignment also are indicated as well as a calculation of the p value.

Full length sequences or fragments of the polynucleotide sequences of the nearest neighbors can be used as probes and primers to identify and isolate the full length sequence of SEQ ID NOs:1-3351. The nearest neighbors can indicate a tissue or cell type to be used to construct a library for the full-length sequences of SEQ ID
5 NOs:1-3351.

EXAMPLE 3

MEMBERS OF PROTEIN FAMILIES

The sequences (SEQ ID NOs:1-3351) were used to conduct a profile
10 search as described in the specification above. Several of the polynucleotides of the invention were found to encode polypeptides having characteristics of a polypeptide belonging to a known protein families (and thus represent new members of these protein families) and/or comprising a known functional domain (Table 5). "Start" and "stop" in Table 3 indicate the position within the individual sequences that align with
15 the query sequence having the indicated SEQ ID NO. The direction indicates the orientation of the query sequence with respect to the individual sequence, where forward (for) indicates that the alignment is in the same direction (left to right) as the sequence provided in the Sequence Listing and reverse (rev) indicates that the alignment is with a sequence complementary to the sequence provided in the Sequence
20 Listing.

Some polynucleotides exhibited multiple profile hits because, for example, the particular sequence contains overlapping profile regions, and/or the sequence contains two different functional domains. These profile hits are described in more detail below.

25 Ank Repeats (ANK). SEQ ID NOs:187, 1268, 1804, 1819, 1830, 1839, 2652, 3015 and 3267 represent polynucleotides encoding an Ank repeat-containing protein. The ankyrin motif is a 33 amino acid sequence named for the protein ankyrin which has 24 tandem 33-amino-acid motifs. Ank repeats were originally identified in the cell-cycle-control protein cdc10 (Breedon et al., *Nature* (1987) 329:651). Proteins
30 containing ankyrin repeats include ankyrin, myotropin, I-kappaB proteins, cell cycle protein cdc10, the Notch receptor (Matsuno et al., *Development* (1997) 124(21):4265); G9a (or BAT8) of the class III region of the major histocompatibility complex (*Biochem J.* 290:811-818, 1993), FABP, GABP, 53BP2, Lin12, glp-1, SW14, and SW16. The functions of the ankyrin repeats are compatible with a role in protein-

protein interactions (Bork, *Proteins* (1993) 17(4):363; Lambert and Bennet, *Eur. J. Biochem.* (1993) 211:1; Kerr et al., *Current Op. Cell Biol.* (1992) 4:496; Bennet et al., *J. Biol. Chem.* (1980) 255:6424).

ATPases Associated with Various Cellular Activities (ATPases).

- 5 Sequences within SEQ ID NOs:431, 639, 2135, 2684, 2859, 3197 and 3266 correspond to a sequence that encodes a novel member of the "ATPases Associated with diverse cellular Activities" (AAA) protein family. The AAA protein family is composed of a large number of ATPases that share a conserved region of about 220 amino acids that contains an ATP-binding site (Froehlich et al., *J. Cell Biol.* (1991) 114:443; Erdmann et al., *Cell* (1991) 64:499; Peters et al., *EMBO J.* (1990) 9:1757; Kunau et al., *Biochimie* 10 (1993) 75:209-224; Confalonieri et al., *BioEssays* (1995) 17:639; <http://yeamob.pci.chemie.uni-tuebingen.de/AAA/Description.html>). The proteins that belong to this family either contain one or two AAA domains. In general, the AAA domains in these proteins act as ATP-dependent protein clamps (Confalonieri et al. 15 (1995) *BioEssays* 17:639). In addition to the ATP-binding 'A' and 'B' motifs, which are located in the N-terminal half of this domain, there is a highly conserved region located in the central part of the domain which was used in the development of the signature pattern. The consensus pattern is: [LIVMT]-x-[LIVMT]-[LIVMF]-x-[GATMC]-[ST]-[NS]-x(4)-[LIVM]-D-x-A-[LIFA]-x-R.

- 20 Bromodomain (bromodomain). SEQ ID NO:1814 represents a polynucleotide encoding a polypeptide having a bromodomain region (Haynes et al., 1992, *Nucleic Acids Res.* 20:2693-2603, Tamkun et al., 1992, *Cell* 68:561-572, and Tamkun, 1995, *Curr. Opin. Genet. Dev.* 5:473-477), which is a conserved region of about 70 amino acids. The bromodomain is thought to be involved in protein-protein 25 interactions and may be important for the assembly or activity of multicomponent complexes involved in transcriptional activation. The consensus pattern, which spans a major part of the bromodomain, is: [STANVF]-x(2)-F-x(4)-[DNS]-x(5,7)-[DENQTF]-Y-[HFY]-x(2)-[LIVMFY]-x(3)-[LIVM]-x(4)-[LIVM]-x(6,8)-Y-x(12,13)-[LIVM]-x(2)-N-[SACF]-x(2)-[FY].

- 30 Basic Region Plus Leucine Zipper Transcription Factors (BZIP). SEQ ID NOs:410, 552, 768, 822, 836, 1288, 1365, 1454, 1540, 1549, 1556, 1557, 1563, 1622, 1630, 1704, 1808, 2363, 2424, 3147, 3152, 3158 and 3208 represent polynucleotides encoding a novel member of the family of basic region plus leucine zipper transcription factors. The bZIP superfamily (Hurst, *Protein Prof.* (1995) 2:105; 35 and Ellenberger, *Curr. Opin. Struct. Biol.* (1994) 4:12) of eukaryotic DNA-binding

transcription factors encompasses proteins that contain a basic region mediating sequence-specific DNA-binding followed by a leucine zipper required for dimerization. The consensus pattern for this protein family is: [KR]-x(1,3)-[RKSAQ]-N-x(2)-[SAQ](2)-x-[RKTAENQ]-x-R-x-[RK].

- 5 EF Hand (EFhand). SEQ ID NOs:820, 1755 and 3285 correspond to polynucleotides encoding a novel protein in the family of EF-hand proteins. Many calcium-binding proteins belong to the same evolutionary family and share a type of calcium-binding domain known as the EF-hand (Kawasaki et al., *Protein. Prof.* (1995) 2:305-490). This type of domain consists of a twelve residue loop flanked on both sides
10 by a twelve residue alpha-helical domain. In an EF-hand loop the calcium ion is coordinated in a pentagonal bipyramidal configuration. The six residues involved in the binding are in positions 1, 3, 5, 7, 9 and 12; these residues are denoted by X, Y, Z, -Y, -X and -Z. The invariant Glu or Asp at position 12 provides two oxygens for liganding Ca (bidentate ligand). The consensus pattern includes the complete EF-hand loop as
15 well as the first residue which follows the loop and which seem to always be hydrophobic: D-x-[DNS]-{ILVFW}-[DENSTG]-[DNQGHRK]-{GP}-[LIVMC]-[DENQSTAGC]-x(2)-[DE]-[LIVMFYW].

- Ets Domain (Ets Nterm). SEQ ID NO:1811 represents a polynucleotide encoding a polypeptide with N-terminal homology in ETS domain. Proteins of this
20 family contain a conserved domain, the "ETS-domain," that is involved in DNA binding. The domain appears to recognize purine-rich sequences; it is about 85 to 90 amino acids in length, and is rich in aromatic and positively charged residues (Wasylyk, et al., *Eur. J. Biochem.* (1993) 211:718). The *ets* gene family encodes a novel class of DNA-binding proteins, each of which binds a specific DNA sequence and comprises an
25 *ets* domain that specifically interacts with sequences containing the common core trinucleotide sequence GGA. In addition to an *ets* domain, native *ets* proteins comprise other sequences which can modulate the biological specificity of the protein. *Ets* genes and proteins are involved in a variety of essential biological processes including cell growth, differentiation and development, and three members are implicated in
30 oncogenic process.

- G-Protein Alpha Subunit (G-alpha). SEQ ID NO:1846 represents a polynucleotide encoding a novel polypeptide of the G-protein alpha subunit family. Guanine nucleotide binding proteins (G-proteins) are a family of membrane-associated proteins that couple extracellularly-activated integral-membrane receptors to
35 intracellular effectors, such as ion channels and enzymes that vary the concentration of

second messenger molecules. G-proteins are composed of 3 subunits (alpha, beta and gamma) which, in the resting state, associate as a trimer at the inner face of the plasma membrane. The alpha subunit binds GTP and exhibits GTPase activity. G-protein alpha subunits are 350-400 amino acids in length and have molecular weights in the range 40-45 kDa. Seventeen distinct types of alpha subunit have been identified in mammals, and fall into 4 main groups on the basis of both sequence similarity and function: alpha-s, alpha-q, alpha-i and alpha-12 (Simon et al., *Science* (1993) 252:802). They are often N-terminally acylated, usually with myristate and/or palmitoylate, and these fatty acid modifications can be important for membrane association and high-affinity interactions with other proteins.

Helicases conserved C-terminal domain (helicase C). SEQ ID NOs:1496, 2826 and 2871 represent polynucleotides encoding novel members of the DEAD/H helicase family. A number of eukaryotic and prokaryotic proteins have been characterized (Schmid S.R., et al., *Mol. Microbiol.* (1992) 6:283; Linder P., et al., *Nature* (1989) 337:121; Wassarman D.A., et al., *Nature* (1991) 349:463) on the basis of their structural similarity. All are involved in ATP-dependent, nucleic-acid unwinding. All DEAD box family members of the above proteins share a number of conserved sequence motifs, some of which are specific to the DEAD family while others are shared by other ATP-binding proteins or by proteins belonging to the helicases 'superfamily' (Hodgman T.C., *Nature* (1988) 333:22 and *Nature* (1988) 333:578 (Errata). One of these motifs, called the "D-E-A-D-box", represents a special version of the B motif of ATP-binding proteins. Some other proteins belong to a subfamily which have His instead of the second Asp and are thus said to be "D-E-A-H-box" proteins (Wassarman D.A., et al., *Nature* (1991) 349:463; Harosh I., et al., *Nucleic Acids Res.* (1991) 19:6331; Koonin E.V. et al., *J. Gen. Virol.* (1992) 73:989. The following signature patterns are used to identify members of both subfamilies: 1) [LIVMF](2)-D-E-A-D-[RKEN]-x-[LIVMFYGSTN]; and 2) [GSAH]-x-[LIVMF](3)-D-E-[ALIV]-H-[NECR].

Homeobox domain (homeobox). SEQ ID NOs:1676, 1820 and 1821 represent polynucleotides encoding proteins having a homeobox domain. The homeobox is a protein domain of 60 amino acids (Gehring In: Guidebook to the Homeobox Genes, Duboule D., Ed., pp. 1-10, Oxford University Press, Oxford, (1994); Buerklin In: Guidebook to the Homeobox Genes, pp25-72, Oxford University Press, Oxford, (1994); Gehring, *Trends Biochem. Sci.* (1992) 17:277-280; Gehring et al., *Annu. Rev. Genet.* (1986) 20:147-173; Schofield, *Trends Neurosci.* (1987) 10:3-6) first

identified in a number of *Drosophila* homeotic and segmentation proteins. It is extremely well conserved in many other animals, including vertebrates. This domain binds DNA through a helix-turn-helix type of structure. Several proteins that contain a homeobox domain play an important role in development. Most of these proteins are sequence-specific DNA-binding transcription factors. The homeobox domain is also very similar to a region of the yeast mating type proteins. These are sequence-specific DNA-binding proteins that act as master switches in yeast differentiation by controlling gene expression in a cell type-specific fashion.

A schematic representation of the homeobox domain is shown below.

The helix-turn-helix region is shown by the symbols 'H' (for helix), and 't' (for turn).

```

XXXXXXXXXXXXXXXXXXXXXXXXXXXXHHHHHHHHHtTtHHHHHHHHHXXXXXXXXX
1                                     60

```

The pattern detects homeobox sequences 24 residues long and spans positions 34 to 57 of the homeobox domain. The consensus pattern is as follows: [LIVMFYGG]-[ASLVR]-x(2)-[LIVMSTACN]-x-[LIVM]-x(4)-[LIV]-[RKNQESTAIY]-[LIVFSTNKH]-W-[FYVC]-x-[NDQTAH]-x(5)-[RKNAIMW].

MAP kinase kinase (mkk). SEQ ID NOs:29, 31, 196, 3175, 3190 and 3281 represent novel members of the MAP kinase kinase family. MAP kinases (MAPK) are involved in signal transduction, and are important in cell cycle and cell growth controls. The MAP kinase kinases (MAPKK) are dual-specificity protein kinases which phosphorylate and activate MAP kinases. MAPKK homologues have been found in yeast, invertebrates, amphibians, and mammals. Moreover, the MAPKK/MAPK phosphorylation switch constitutes a basic module activated in distinct pathways in yeast and in vertebrates. MAPKKs are essential transducers through which signals must pass before reaching the nucleus. For review, see, *e.g.*, Biologie *Biol Cell* (1993) 79:193-207; Nishida et al., *Trends Biochem Sci* (1993) 18:128-31; Ruderman, *Curr Opin Cell Biol* (1993) 5:207-13; Dhanasekaran et al., *Oncogene* (1998) 17:1447-55; Kiefer et al., *Biochem Soc Trans* (1997) 25:491-8; and Hill, *Cell Signal* (1996) 8:533-44.

Protein Kinase (protkinase). SEQ ID NOs:1157, 1478, 1496, 2286, 2969 and 3190 represent polynucleotides encoding protein kinases. Protein kinases catalyze phosphorylation of proteins in a variety of pathways, and are implicated in cancer. Eukaryotic protein kinases (Hanks S.K., et al., *FASEB J.* (1995) 9:576; Hunter T., *Meth. Enzymol.* (1991) 200:3; Hanks S.K., et al., *Meth. Enzymol.* (1991) 200:38; Hanks S.K.,

- Curr. Opin. Struct. Biol.* (1991) 1:369; Hanks S.K. et al., *Science* (1988) 241:42) are enzymes that belong to a very extensive family of proteins which share a conserved catalytic core common to both serine/threonine and tyrosine protein kinases. There are a number of conserved regions in the catalytic domain of protein kinases. The first region, which is located in the N-terminal extremity of the catalytic domain, is a glycine-rich stretch of residues in the vicinity of a lysine residue, which has been shown to be involved in ATP binding. The second region, which is located in the central part of the catalytic domain, contains a conserved aspartic acid residue which is important for the catalytic activity of the enzyme (Knighton D.R. et al., *Science* (1991) 253:407).
- 5 The protein kinase profile includes two signature patterns for this second region: one specific for serine/threonine kinases and the other for tyrosine kinases. A third profile is based on the alignment in (Hanks S.K. et al., *FASEB J.* (1995) 9:576) and covers the entire catalytic domain.

- The consensus patterns are as follows: 1) [LIV]-G-{P}-G-{P}-[FYWMGSTNH]-[SGA]-{PW}-[LIVCAT]-{PD}-x-[GSTACLIVMFY]-x(5,18)-[LIVMFYWCSTAR]-[AIVP]-[LIVMFAGCKR]-K, where K binds ATP; 2) [LIVMFYCY]-x-[HY]-x-D-[LIVMFY]-K-x(2)-N-[LIVMFYCT](3), where D is an active site residue; and 3) [LIVMFYCY]-x-[HY]-x-D-[LIVMFY]-[RSTAC]-x(2)-N-[LIVMFYCY], where D is an active site residue.

- 20 If a protein analyzed includes two of the above protein kinase signatures, the probability of it being a protein kinase is close to 100%.

- Ras family proteins (ras). SEQ ID NOs:1688 and 3258 represent polynucleotides encoding novel members of the ras family of small GTP/GDP-binding proteins (Valencia et al., 1991, *Biochemistry* 30:4637-4648). Ras family members generally require a specific guanine nucleotide exchange factor (GEF) and a specific GTPase activating protein (GAP) as stimulators of overall GTPase activity. Among ras-related proteins, the highest degree of sequence conservation is found in four regions that are directly involved in guanine nucleotide binding. The first two constitute most of the phosphate and Mg²⁺ binding site (PM site) and are located in the first half of the G-domain. The other two regions are involved in guanosine binding and are located in the C-terminal half of the molecule. Motifs and conserved structural features of the ras-related proteins are described in Valencia et al., 1991, *Biochemistry* 30:4637-4648. A major consensus pattern of ras proteins is: D-T-A-G-Q-E-K-[LF]-G-G-L-R-[DE]-G-Y-Y.

Thioredoxin family active site (Thioredoxin). SEQ ID NO:1677 represents a polynucleotide encoding a protein having a thioredoxin family active site. Thioredoxins (Holmgren A., *Annu. Rev. Biochem.* (1985) 54:237; Gleason F.K. et al., *FEMS Microbiol. Rev.* (1988) 54:271; Holmgren, A. *J. Biol. Chem.* (1989) 264:13963; 5 Eklund H. et al., *Proteins* (1991) 11:13) are small proteins of approximately one hundred amino- acid residues which participate in various redox reactions via the reversible oxidation of an active center disulfide bond. They exist in either a reduced form or an oxidized form where the two cysteine residues are linked in an intramolecular disulfide bond. Thioredoxin is present in prokaryotes and eukaryotes 10 and the sequence around the redox-active disulfide bond is well conserved. All PDI contains two or three (ERp72) copies of the thioredoxin domain. The consensus pattern is: [LIVMF]-[LIVMSTA]-x-[LIVMFYC]-[FYWSTHE]-x(2)-[FYWGNT]-C-[GATPLVE]-[PHYWSTA]-C-x(6)-[LIVMFYWT] (where the two C's form the redox-active bond).

15 Trypsin (trypsin). SEQ ID NO:1410 corresponds to a novel serine protease of the trypsin family. The catalytic activity of the serine proteases from the trypsin family is provided by a charge relay system involving an aspartic acid residue hydrogen-bonded to a histidine, which itself is hydrogen-bonded to a serine. The sequences in the vicinity of the active site serine and histidine residues are well 20 conserved in this family of proteases (Brenner S., *Nature* (1988) 334:528). The consensus patterns for this trypsin protein family are: 1) [LIVM]-[ST]-A-[STAG]-H-C, where H is the active site residue; and 2) [DNSTAGC]-[GSTAPIMVQH]-x(2)-G-[DE]-S-G-[GS]-[SAPHV]- [LIVMFYWH]-[LIVMFYSTANQH], where S is the active site residue. All sequences known to belong to this family are detected by the above 25 consensus sequences, except for 18 different proteases which have lost the first conserved glycine. If a protein includes both the serine and the histidine active site signatures, the probability of it being a trypsin family serine protease is 100%.

WD Domain, G-Beta Repeats (WD domain). SEQ ID NOs:1336, 1380, 1711, 1762, 1909, 2218, 3047, 3108 and 3292 represent novel members of the WD 30 domain/G-beta repeat family. Beta-transducin (G-beta) is one of the three subunits (alpha, beta, and gamma) of the guanine nucleotide-binding proteins (G proteins) which act as intermediaries in the transduction of signals generated by transmembrane receptors (Gilman, *Annu. Rev. Biochem.* (1987) 56:615). The alpha subunit binds to and hydrolyzes GTP; the functions of the beta and gamma subunits are less clear but 35 they seem to be required for the replacement of GDP by GTP as well as for membrane

anchoring and receptor recognition. In higher eukaryotes, G-beta exists as a small multigene family of highly conserved proteins of about 340 amino acid residues. Structurally, G-beta consists of eight tandem repeats of about 40 residues, each containing a central Trp-Asp motif (this type of repeat is sometimes called a WD-40 repeat). The consensus pattern for the WD domain/G-Beta repeat family is:
 5 [LIVMSTAC]-[LIVMFYWSTAGC]-[LIMSTAG]-[LIVMSTAGC]-x(2)-[DN]-x(2)-
 [LIVMWSTAC]-x-[LIVMFSTAG]-W-[DEN]-[LIVMFSTAGCN].

wnt Family of Developmental Signaling Proteins (Wnt dev sign). SEQ ID NO:1538 corresponds to a novel member of the wnt family of developmental signaling proteins. Wnt-1 (previously known as int-1), the seminal member of this family, (Nusse R., *Trends Genet.* (1988) 4:291) is thought to play a role in intercellular communication and seems to be a signalling molecule important in the development of the central nervous system (CNS). All wnt family proteins share the following features characteristics of secretory proteins: a signal peptide, several potential N-glycosylation sites and 22 conserved cysteines that are probably involved in disulfide bonds. The Wnt proteins seem to adhere to the plasma membrane of the secreting cells and are therefore likely to signal over only few cell diameters. The consensus pattern, which is based upon a highly conserved region including three cysteines, is as follows: C-K-C-H-G-[LIVMT]-S-G-x-C.
 15

Protein Tyrosine Phosphatase (Y phosphatase). SEQ ID NO:1417 represents a polynucleotide encoding a protein tyrosine kinase. Tyrosine specific protein phosphatases (EC 3.1.3.48) (PTPase) (Fischer et al., *Science* (1991) 253:401; Charbonneau et al., *Annu. Rev. Cell Biol.* (1992) 8:463; Trowbridge, *J. Biol. Chem.* (1991) 266:23517; Tonks et al., *Trends Biochem. Sci.* (1989) 14:497; and Hunter, *Cell* (1989) 58:1013) catalyze the removal of a phosphate group attached to a tyrosine residue. These enzymes are very important in the control of cell growth, proliferation, differentiation and transformation. Multiple forms of PTPase have been characterized and can be classified into two categories: soluble PTPases and transmembrane receptor proteins that contain PTPase domain(s). Structurally, all known receptor PTPases are made up of a variable length extracellular domain, followed by a transmembrane region and a C-terminal catalytic cytoplasmic domain. PTPase domains consist of about 300 amino acids. The search of two conserved cysteines has been shown to be absolutely required for activity. Furthermore, a number of conserved residues in its immediate vicinity have also been shown to be important. The consensus pattern for PTPases is:
 25 [LIVMF]-H-C-x(2)-G-x(3)-[STC]-[STAGP]-x-[LIVMFY]; C is the active site residue.
 30
 35

Zinc Finger, C2H2 Type (Zincfing C2H2). SEQ ID NOs:308, 807, 1324, 1503, 1527, 3081, 3193 and 3306 correspond to polynucleotides encoding novel members of the of the C2H2 type zinc finger protein family. Zinc finger domains (Klug et al., *Trends Biochem. Sci.* (1987) 12:464; Evans et al., *Cell* (1988) 52:1; Payre et al., 5 *FEBS Lett.* (1988) 234:245; Miller et al., *EMBO J.* (1985) 4:1609; and Berg, *Proc. Natl. Acad. Sci. USA* (1988) 85:99) are nucleic acid-binding protein structures. In addition to the conserved zinc ligand residues, it has been shown that a number of other positions are also important for the structural integrity of the C2H2 zinc fingers. (Rosenfeld et al., *J. Biomol. Struct. Dyn.* (1993) 11:557) The best conserved position is found four 10 residues after the second cysteine; it is generally an aromatic or aliphatic residue. The consensus pattern for C2H2 zinc fingers is: C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H. The two C's and two H's are zinc ligands.

Src homology 2. SEQ ID NOs:186, 2591, 3307 and 3339 represent polynucleotides encoding novel members of the family of Src homology 2 (SH2) 15 proteins. The Src homology 2 (SH2) domain is a protein domain of about 100 amino acid residues first identified as a conserved sequence region between the oncoproteins Src and Fps (Sadowski I. et al., *Mol. Cell. Biol.* 6:4396-4408 (1986)). Similar sequences are found in many other intracellular signal-transducing proteins (Russel R.B. et al., *FEBS Lett.* 304:15-20 (1992)). SH2 domains function as regulatory modules of 20 intracellular signalling cascades by interacting with high affinity to phosphotyrosine-containing target peptides in a sequence-specific and phosphorylation-dependent manner (Marangere L.E.M., Pawson T., *J. Cell Sci. Suppl.* 18:97-104 (1994); Pawson T., Schlessinger J., *Curr. Biol.* 3:434-442 (1993); Mayer B.J., Baltimore D., *Trends Cell. Biol.* 3:8-13 (1993); Pawson T., *Nature* 373:573-580 (1995)).

25 The SH2 domain has a conserved 3D structure consisting of two alpha helices and six to seven beta-strands. The core of the domain is formed by a continuous beta-meander composed of two connected beta-sheets (Kuriyan J., Cowburn D., *Curr. Opin. Struct. Biol.* 3:828-837(1993)). The profile to detect SH2 domains is based on a structural alignment consisting of 8 gap-free blocks and 7 linker regions totaling 92 30 match positions.

Src homology 3. SEQ ID NO:234, 1832, and 1835 represent polynucleotides encoding novel members of the family of Src homology 3 (SH3) proteins. The Src homology 3 (SH3) domain is a small protein domain of about 60 amino acid residues first identified as a conserved sequence in the non-catalytic part of 35 several cytoplasmic protein tyrosine kinases (e.g., Src, Abl, Lck) (Mayer B.J. et al.,

Nature 332:272-275 (1988)). Since then, it has been found in a great variety of other intracellular or membrane-associated proteins (Musacchio A. et al., *FEBS Lett.* 307:55-61 (1992); Pawson T., Schlessinger J., *Curr. Biol.* 3:434-442 (1993); Mayer B.J., Baltimore D., *Trends Cell Biol.* 3:8-13 (1993); Pawson T., *Nature* 373:573-580 (1995)).

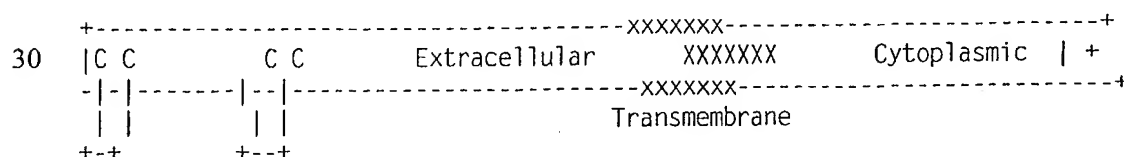
5 The SH3 domain has a characteristic fold which consists of five or six beta strands arranged as two tightly packed anti-parallel beta sheets. The linker regions may contain short helices (Kuriyan J., Cowburn D., *Curr. Opin. Struct. Biol.* 3:828-837 (1993)).

The function of the SH3 domain may be to mediate assembly of specific
10 protein complexes via binding to proline-rich peptides (Morton C.J., Campbell I.D.,
Curr. Biol. 4:615-617 (1994)).

In general SH3 domains are found as single copies in a given protein, but there are a significant number of proteins with two SH3 domains and a few with 3 or 4 copies.

15 Fibronectin type III. SEQ ID NOs:746 and 1192 represent
polynucleotides encoding novel members of the family of fibronectin type III proteins.
A number of receptors for lymphokines, hematopoietic growth factors and growth
hormone-related molecules have been found to share a common binding domain.
(Bazan J.F., *Biochem. Biophys. Res. Commun.* 164:788-795 (1989); Bazan J.F., *Proc.*
20 *Natl. Acad. Sci. U.S.A.* 87:6934-6938 (1990); Cosman D. et al., *Trends Biochem. Sci.*
15:265-270 (1990); d'Andrea A.D., Fasman G.D., Lodish H.F., *Cell* 58:1023-1024
(1989); d'Andrea A.D., Fasman G.D., Lodish H.F., *Curr. Opin. Cell Biol.* 2:648-651
(1990)).

The conserved region constitutes all or part of the extracellular ligand-binding region and is about 200 amino acid residues long. In the N-terminal of this domain there are two pairs of cysteines known, in the growth hormone receptor, to be involved in disulfide bonds.



Two patterns detect this family of receptors. The first one is derived from the first N-terminal disulfide loop, the second is a tryptophan-rich pattern located at the C-terminal extremity of the extracellular region.

63

A consensus for this protein family is: C-[LVFYR]-x(7,8)-[STIVDN]-C-x-W (The two C's are linked by a disulfide bond]. A second consensus for this protein family is: [STGL]-x-W-[SG]-x-W-S.

LIM domain containing proteins. SEQ ID NOs:1269, 1309, 1360, and 1386 represent polynucleotides encoding novel members of the family of LIM domain containing proteins. A number of proteins contain a conserved cysteine-rich domain of about 60 amino-acid residues. (Freyd G. et al., *Nature* 344:876-879 (1990); Baltz R. et al., *Plant Cell* 4:1465-1466 (1992); Sanchez-Garcia I., Rabbitts T.H., *Trends Genet.* 10:315-320 (1994)).

In the LIM domain, there are seven conserved cysteine residues and a histidine. The arrangement followed by these conserved residues is C-x(2)-C-x(16,23)-H-x(2)-[CH]-x(2)-C-x(2)-C-x(16,21)-C-x(2,3)-[CHD]. The LIM domain binds two zinc ions (Michelsen J.W. et al., *Proc. Natl. Acad. Sci. U.S.A.* 90:4404-4408 (1993)). LIM does not bind DNA, rather it seems to act as interface for protein-protein interaction. The consensus for this protein family is: C-x(2)-C-x(15,21)-[FYWH]-H-x(2)-[CH]-x(2)-C-x(2)-C-x(3)-[LIVMF]. The 5 C's and the H bind zinc.

C2 domain (protein kinase C like). SEQ ID NOs:1325 and 2282 represent polynucleotides encoding novel members of the family of C2 domain containing proteins. Some isozymes of protein kinase C (PKC) contain a domain, known as C2, of about 116 amino-acid residues, which is located between the two copies of the C1 domain (that bind phorbol esters and diacylglycerol) and the protein kinase catalytic domain. (Azzi A. et al., *Eur. J. Biochem.* 208:547-557 (1992); Stabel S., *Semin. Cancer Biol.* 5:277-284 (1994)).

The C2 domain is involved in calcium-dependent phospholipid binding (Davletov B.A., Suedhof T.C., *J. Biol. Chem.* 268:26386-26390 (1993)). Since domains related to the C2 domain are also found in proteins that do not bind calcium, other putative functions for the C2 domain include binding to inositol-1,3,5-tetraphosphate. (Fukuda M., et al., *J. Biol. Chem.* 269:29206-29211 (1994)).

The consensus pattern for the C2 domain is located in a conserved part of that domain, the connecting loop between beta strands 2 and 3. The profile for the C2 domain covers the total domain. The consensus for this protein family is:: [ACG]-x(2)-L-x(2,3)-D-x(1,2)-[NGSTLIF]-[GTMR]-x-[STAP]-D-[PA]-[FY]

Serine proteases, trypsin family, active sites. SEQ ID NO:1410 represents a polynucleotide encoding a novel member of the family of serine protease, trypsin proteins. The catalytic activity of the serine proteases from the trypsin family is

provided by a charge relay system involving an aspartic acid residue hydrogen-bonded to a histidine, which itself is hydrogen-bonded to a serine. The sequences in the vicinity of the active site serine and histidine residues are well conserved in this family of proteases (Brenner S., *Nature* 334:528-530 (1988)).

- 5 A consensus for this protein family is: [LIVM]-[ST]-A-[STAG]-H-C [H is the active site residue]. A second consensus for this protein family is: [DNSTAGC]-[GSTAPIMVQH]-x(2)-G-[DE]-S-G-[GS]-[SAPHV]-[LIVMFYWH]-[LIVMFYSTANQH] [S is the active site residue].

RNA Recognition Motif Domain (RRM, RBD, or RNP). SEQ ID NOS:

- 10 1464 and 1514 represent polynucleotides encoding novel members of the family of RNA recognition motif domain proteins (Bandziulis R.J. et al., *Genes Dev.* 3:431-437 (1989); Dreyfuss G. et al., *Trends Biochem. Sci.* 13:86-91 (1988)).

- Inside the putative RNA-binding domain there are two regions which are highly conserved. The first one is a hydrophobic segment of six residues (which is called the RNP-2 motif); the second one is an octapeptide motif (which is called RNP-1 or RNP-CS). The position of both motifs in the domain is shown in the following schematic representation:
- 15

20 xxxxxxxx#####xx#####xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx
RNP-2 RNP-1

As a consensus pattern for this type of domain the RNP-1 motif was used. The consensus for this protein family is: [RK]-G-{EDRKHPCG}-[AGSCI]-[FY]-[LIVA]-x-[FYLM]

- 25 Phosphatidylinositol-specific phospholipase C, Y Domain. SEQ ID NO: 1707 represents a polynucleotide encoding a novel member of the phosphatidylinositol-specific phospholipase C, Y domain family of proteins. Phosphatidylinositol-specific phospholipase C (EC3.1.4.11), a eukaryotic intracellular enzyme, plays an important role in signal transduction processes (Meldrum E. et al., *Biochim. Biophys. Acta* 1092:49-71 (1991)). It catalyzes the hydrolysis of 1-phosphatidyl-D-myo-inositol-3,4,5- triphosphate into the second messenger molecules diacylglycerol and inositol-1,4,5-triphosphate. This catalytic process is tightly regulated by reversible phosphorylation and binding of regulatory proteins (Rhee S.G., Choi K.D., *Adv. Second Messenger Phosphoprotein Res.* 26:35-61 (1992); Rhee S.G., Choi K.D., *J. Biol. Chem.* 267:12393-12396 (1992); Sternweis P.C., Smrcka A.V., *Trends Biochem. Sci.* 17:502-506 (1992)).
- 30
- 35

All eukaryotic PI-PLCs contain two regions of homology, referred to as "X-box" and "Y-box". The order of these two regions is the same (NH₂-X-Y-COOH), but the spacing is variable. In most isoforms, the distance between these two regions is only 50-100 residues but in the gamma isoforms one PH domain, two SH2 domains, and one SH3 domain are inserted between the two PLC-specific domains. The two conserved regions have been shown to be important for the catalytic activity. At the C-terminal of the Y-box, there is a C2 domain possibly involved in Ca-dependent membrane attachment.

Serine Carboxypeptidases. SEQ ID NO:1744 represents a polynucleotide encoding a novel member of the serine carboxypeptidases family of proteins. Carboxypeptidases may be either metallo carboxypeptidases or serine carboxypeptidases (EC 3.4.16.5 and EC 3.4.16.6). The catalytic activity of the serine carboxypeptidases, like that of the trypsin family serine proteases, is provided by a charge relay system involving an aspartic acid residue hydrogen-bonded to a histidine, which is itself hydrogen-bonded to a serine (Liao D.I., Remington S.J., *J. Biol. Chem.* 265:6528-6531 (1990)).

The sequences surrounding the active site serine and histidine residues are highly conserved in all these serine carboxypeptidases. A consensus for this protein family is: [LIVM]-x-[GTA]-E-S-Y-[AG]-[GS] [S is the active site residue]. A second consensus for this protein family is: [LIVF]-x(2)-[LIVSTA]-x-[IVPST]-x-[GSDNQL]-[SAGV]-[SG]-H-x- [IVAQ]-P-x(3)-[PSA] [H is the active site residue].

dsrm Double-Stranded RNA Binding Motif. SEQ ID NO:1818 represents a polynucleotide encoding a novel member of the dsrm double-stranded RNA binding motif proteins. In eukaryotic cells, a multitude of RNA-binding proteins play key roles in the posttranscriptional regulation of gene expression. Characterization of these proteins has led to the identification of several RNA-binding motifs. Several human and other vertebrate genetic disorders are caused by aberrant expression of RNA-binding proteins. (C. G. Burd & G. Dreyfuss, *Science* 265: 615-621 (1994)).

Proteins containing double stranded RNA binding motifs bind to specific RNA targets. Double stranded RNA binding motifs are exemplified by interferon-induced protein kinase in humans, which is part of the cellular response to dsRNA.

SEQ ID NOs:2577, 3183 and 3195 encode members of the 4 trans-membrane integral membrane protein family. This family consists of type III proteins, which are integral membrane proteins that contain a N-terminal membrane-anchoring domain that is not cleaved during biosynthesis, and which functions as a translocation

signal and a membrane anchor. The proteins also have three additional transmembrane regions. The consensus pattern is: G-x(3)-[LIVMF]-x(2)-[GSA]-[LIVMF] (2)-G-C-x-[GA]-[STA]-x(20-[eG]-x(20-[CwN]-[LIVM](2).

SEQ ID NO:2944 encodes a polypeptide having a calpain large subunit, domain III. Calpains are a family of intracellular proteases that play a variety of biological roles. Calpain 3, also known as p94, is predominantly expressed in skeletal muscle and plays a role in limb-girdle muscular dystrophy type 2A. (Sorimachi, H. et al., *Biochem. J.* 328:721-732, 1997).

SEQ ID NOs:1911 and 1980 encode polypeptides having a C3HC4 type zinc finger domain (RING finger), which is a cysteine-rich domain of 40 to 60 residues that binds two atoms of zinc, and is believed to be involved in mediating protein-protein interactions. Mammalian proteins of this family include V(D)J recombination activating protein, which activates the rearrangement of immunoglobulin and T-cell receptor genes; breast cancer type 1 susceptibility protein (BRCA1); bmi-1 proto-oncogene; cbl proto-oncogene; and mel-18 protein, which is expressed in a variety of tumor cells and is a transcriptional repressor that recognizes and binds a specific DNA sequence. The consensus pattern is: C-x-H-x-[LIVMFY]-C-x(2)-C-[LIVMYA].

SEQ ID NO:3274 encodes a eukaryotic transcription factor with a fork head domain, of about 100 amino acid residues. Proteins of this group are transcription factors, including mammalian transcription factors HNF-3-alpha, -beta, and -gamma; interleukin-enhancer binding factor; and HTLF, which binds to a region of human T-cell leukemia virus long terminal repeat. The consensus pattern is [KR]-P-[PTQ]-[FYLVQH]-S-[FY]x(2)-[LIVM]-X(3,4)-[AC]-[LIM].

SEQ ID NO:3345 encodes a polypeptide having a PDZ domain. Several dozen signaling proteins belong to this group of proteins that have 80-100 residue repeats known as PDZ domains. Several of the proteins interact with the C-terminal tetrapeptide motifs X-Ser/Thr/X-Val-COO- of ion channels and/or receptors. (Ponting, C. P., *Protein Sci.* 6:464-468, 1997.)

SEQ ID NO:3351 encodes a polypeptide in the family of phorbol esters/glycerol binding proteins. Phorbol esters (PE) are analogues of diacylglycerol (DAG) and potent tumor promoters. DAG activates a family of serine-threonine protein kinases, known as protein kinase C. The N-terminal region of protein kinase C binds PE and DAG, and contains one or two copies of a cysteine-rich domain of about 50 amino acid residues. Other proteins having this domain include diacylglycerol kinase; the vav oncogene; and N-chimaerin, a brain-specific protein. The DAG/PE binding

domain binds two zinc ions through the six cysteines and two histidines that are conserved in the domain. The consensus pattern is: H-x-[LIVMFYW]-x(8, 11)-C-x(2)-C-x(3)-[LIVMFC]-x(5, 10)-C-x(2)-C-x(4)-[HD]-x(2)-C-x(5, 9)-C.

5 SEQ ID NO:2216 encodes a polypeptide having a WW/rsp5/WWP domain. The protein is named for the presence of conserved aromatic positions, generally tryptophan, as well as a conserved proline. Proteins having the domain include dystrophin, vertebrate YAP protein, and IQGAP, a human GTPase activating protein which acts on ras. The consensus pattern is: W-x(9,11)-[VFY]-[FYW]-x(6,7)-[GSTNE]-[GSTQCR]-[FYW]-x(2)-P.

10 SEQ ID NO:2428 encodes a member of the dual specificity phosphatase family, having a catalytic domain, and SEQ IDS NOs:2281 and 2310 encode members of the protein tyrosine phosphatase family. These families are related and classified as tyrosine specific protein phosphatases. The enzymes catalyze the removal of a phosphate group from a tyrosine residue, and are important in the control of cell growth,
15 proliferation, differentiation, and transformation. The consensus pattern is [LIVMF]-H-C-x(2)-G-x(3)-[STC]-[STAGP]-x-[LIVMFY].

| SEQ ID | CLUSTER | SEQ NAME | ORIENTATION | CLONE ID | LIBRARY |
|--------|---------|---------------------------|-------------|----------------|---------|
| 565 | 374502 | RTA00002673F.i.08.2.P.Seq | F | M00039080C:H06 | CH09LNL |
| 566 | 240615 | RTA00002672F.e.19.2.P.Seq | F | M00038995D:E05 | CH09LNL |
| 567 | 379207 | RTA00002670F.b.07.2.P.Seq | F | M00033306D:G08 | CH09LNL |
| 568 | 427893 | RTA00002665F.k.19.1.P.Seq | F | M00031419D:C04 | CH08LNL |
| 569 | 377530 | RTA00002684F.g.19.2.P.Seq | F | M00040305A:D11 | CH09LNL |
| 570 | 429707 | RTA00002668F.c.11.1.P.Seq | F | M00032918C:B10 | CH08LNL |
| 571 | 427610 | RTA00002665F.i.04.1.P.Seq | F | M00028770A:D04 | CH08LNL |
| 572 | 100699 | RTA00002662F.b.22.2.P.Seq | F | M00006680B:D02 | CH02COH |
| 573 | 378974 | RTA00002682F.m.21.1.P.Seq | F | M00040017A:C06 | CH09LNL |
| 574 | 373607 | RTA00002674F.d.15.1.P.Seq | F | M00039127D:E10 | CH09LNL |
| 575 | 262951 | RTA00002665F.d.04.3.P.Seq | F | M00028215D:F03 | CH08LNL |
| 576 | 30748 | RTA00002713F.e.11.1.P.Seq | F | M00027301B:B08 | CH04MAL |
| 577 | 161116 | RTA00002714F.c.11.1.P.Seq | F | M00027837C:D09 | CH04MAL |
| 578 | 379211 | RTA00002682F.p.20.1.P.Seq | F | M00040029A:G04 | CH09LNL |
| 579 | 430689 | RTA00002669F.i.24.1.P.Seq | F | M00033243B:A05 | CH08LNL |
| 580 | 374122 | RTA00002673F.l.22.2.P.Seq | F | M00039104D:C09 | CH09LNL |
| 581 | 376521 | RTA00002677F.h.06.2.P.Seq | F | M00039398A:B10 | CH09LNL |
| 582 | 372834 | RTA00002670F.b.12.2.P.Seq | F | M00033308B:G05 | CH09LNL |
| 583 | 379014 | RTA00002682F.o.02.1.P.Seq | F | M00040022C:D06 | CH09LNL |
| 584 | 376344 | RTA00002677F.b.18.2.P.Seq | F | M00039340B:G08 | CH09LNL |
| 585 | 376485 | RTA00002676F.f.01.2.P.Seq | F | M00039288C:B11 | CH09LNL |
| 586 | 21661 | RTA00002709F.e.18.1.P.Seq | F | M00005820C:E04 | CH02COH |
| 587 | 376539 | RTA00002675F.b.18.1.P.Seq | F | M00039211A:C12 | CH09LNL |
| 588 | 431645 | RTA00002669F.h.15.3.P.Seq | F | M00035223B:H07 | CH08LNL |
| 589 | 163293 | RTA00002714F.c.20.1.P.Seq | F | M00028120D:F12 | CH04MAL |
| 590 | 178614 | RTA00002713F.c.20.1.P.Seq | F | M00027263A:F10 | CH04MAL |
| 591 | 373274 | RTA00002670F.i.22.2.P.Seq | F | M00033432B:H10 | CH09LNL |
| 592 | 379820 | RTA00002679F.f.15.1.P.Seq | F | M00039677A:B08 | CH09LNL |
| 593 | 160536 | RTA00002663F.f.10.1.P.Seq | F | M00022233C:A12 | CH03MAH |
| 594 | 373313 | RTA00002671F.m.02.1.P.Seq | F | M00038328D:A05 | CH09LNL |
| 595 | 26429 | RTA00002712F.k.23.1.P.Seq | F | M00027022D:G11 | CH04MAL |
| 596 | 17983 | RTA00002711F.f.10.1.P.Seq | F | M00022979A:D05 | CH03MAH |
| 597 | 375388 | RTA00002681F.j.22.2.P.Seq | F | M00039888B:D03 | CH09LNL |
| 598 | 63005 | RTA00002712F.m.21.1.P.Seq | F | M00027094A:B03 | CH04MAL |
| 599 | 23030 | RTA00002709F.b.10.2.P.Seq | F | M00005384A:C11 | CH02COH |
| 600 | 372946 | RTA00002670F.l.07.2.P.Seq | F | M00033457D:A05 | CH09LNL |
| 601 | 375351 | RTA00002680F.e.15.2.P.Seq | F | M00039792A:B04 | CH09LNL |
| 602 | 374502 | RTA00002673F.i.08.1.P.Seq | F | M00039080C:H06 | CH09LNL |
| 603 | 376911 | RTA00002682F.e.09.1.P.Seq | F | M00039938C:A08 | CH09LNL |
| 604 | 376024 | RTA00002675F.n.15.1.P.Seq | F | M00039257D:C03 | CH09LNL |
| 605 | 377194 | RTA00002679F.h.20.1.P.Seq | F | M00039685A:A08 | CH09LNL |
| 606 | 379643 | RTA00002682F.g.08.1.P.Seq | F | M00039978A:G03 | CH09LNL |
| 607 | 379610 | RTA00002680F.k.11.1.P.Seq | F | M00039815C:F09 | CH09LNL |
| 608 | 25613 | RTA00002711F.g.06.1.P.Seq | F | M00023024D:F12 | CH03MAH |
| 609 | 207466 | RTA00002664F.j.08.1.P.Seq | F | M00027733A:A02 | CH04MAL |
| 610 | 400052 | RTA00002687F.h.13.1.P.Seq | F | M00040291D:C05 | CH14EDT |
| 611 | 21290 | RTA00002712F.g.01.1.P.Seq | F | M00026859D:D01 | CH04MAL |

| SEQ ID | Nearest Neighbor (BlastN vs. Genbank) | | | Nearest Neighbor (BlastX vs. Non-Redundant Proteins) | | |
|--------|---------------------------------------|--|---------|--|--|---------|
| | ACCESSION | DESCRIPTION | P VALUE | ACCESSION | DESCRIPTION | P VALUE |
| | | Human WD protein | | | mucin - human >gi 501033 | |
| 573 | U57058 | IR10 pre-mRNA, partial cds | 0.19 | 631302 | (U14383) mucin [Homo sapiens] | 0.60 |
| 574 | AF034460 | Penicillium thomii internal transcribed spacer 1, 5.8S ribosomal RNA gene and internal transcribed spacer 2, complete sequence; and 28S ribosomal RNA gene, partial sequence | 0.19 | 114136 | AMINO-ACID ACETYLTRANSFERASE Pseudomonas aeruginosa >gi 151036 (M38358) N-acetylglutamate synthase [Pseudomonas aeruginosa] | 0.35 |
| 575 | U95098 | Xenopus laevis mitotic phosphoprotein 44 mRNA, partial cds | 0.19 | 105270 | alpha-2-adrenergic receptor - human name 'ADRA2R' [Homo sapiens] | 0.27 |
| 576 | AG001475 | Homo sapiens genomic DNA, 21q region, clone: 125H6N2 | 0.19 | 94977 | hypothetical protein 3 - Pseudomonas sp. (DSM 6898) plasmid pKB740 >gi 45867 (X66604) ORF3 | 0.16 |
| 577 | M63284 | Mouse IgG receptor (beta-Fc-gamma-RII) gene, exons 9 and 10, clones lambda-Fc(3.2.93). | 0.19 | 3024681 | TRANSCRIPTION INITIATION FACTOR TFIID 135 KD SUBUNIT (TAFII-135) (TAFII135) (TAFII-130) of RNA polymerase II transcription factor TFIID [Homo sapiens] | 0.088 |
| 578 | U38241 | Pseudomonas aeruginosa orotate phosphoribosyl transferase (pyrE), catabolite repression control protein (crc) and RNasePH (rph) genes, complete cds | 0.19 | 3044086 | (AF055904) unknown [Myxococcus xanthus] | 0.052 |
| 579 | AF039734 | Lontra longicaudis transthyretin intron 1, partial sequence | 0.19 | 322759 | pistil extensin-like protein (clone pMG14) - common tobacco (fragment) >gi 19927 (Z14015) pistil extensin like protein [Nicotiana tabacum] | 0.030 |
| 580 | U95094 | Xenopus laevis XL-INCENP (XL-INCENP) mRNA, complete cds | 0.19 | 2147194 | collagen - Paralvinella grasslei | 0.002 |

| SEQ ID | Nearest Neighbor (BlastN vs. Genbank) | | | Nearest Neighbor (BlastX vs. Non-Redundant Proteins) | | |
|--------|---------------------------------------|--|---------|--|---|---------|
| | ACCESSION | DESCRIPTION | P VALUE | ACCESSION | DESCRIPTION | P VALUE |
| 576 | AB002333 | Human mRNA for KIAA0335 gene, complete cds | 0.19 | <NONE> | <NONE> | <NONE> |
| 577 | U53566 | Macaca mulatta pit-1/GHF-1 transcription factor mRNA, complete cds | 0.19 | 1078068 | probable membrane protein YLR311c - yeast | 9.2 |
| 578 | U73664 | Human t(11;14)(q13;q32) breakpoint junction sequence | 0.19 | 116734 | COAT PROTEIN (CAPSID PROTEIN) virus >gi 58901 (X62133) CyMV coat protein gene product | 8.8 |
| 579 | AF004054 | Heterophyllaea pustulata rps16 gene, chloroplast gene, partial intron sequence | 0.19 | 1928991 | (U92815) heat shock protein 70 precursor [Citrullus lanatus] | 8.7 |
| 580 | Z27081 | Caenorhabditis elegans cosmid M01A8, complete sequence [Caenorhabditis elegans] | 0.19 | 2496247 | HYPOTHETICAL ATP-BINDING PROTEIN MJ0625 >gi 2128413 pir A64378 hypothetical protein MJ0625 - Methanococcus jannaschii >gi 1591336 (U67510) M. jannaschii predicted coding region MJ0625 | 8.6 |
| 581 | Z74145 | S.cerevisiae chromosome IV reading frame ORF YDL097c | 0.19 | 1174425 | TYROSINE-PROTEIN KINASE SPK-1 | 6.7 |
| 582 | D38547 | Small round structured virus genomic RNA, 3'terminal sequence containing ORF2 and ORF3 | 0.19 | 971318 | (Z48053) putative protein [Bovine herpesvirus 1] | 5.1 |

| | | | | | | |
|-------------|-------------|-------------|-------------------|-------------|-------------|-----|
| ggaggggagg | ggaggggggtg | gcacccctggc | ctctaggata | aatgcctgga | gtatagggca | 240 |
| gcgcccacggg | cacttggaga | ccctgtcctg | cgcactctgcc | aagcctggca | gtttttagag | 300 |
| ttttttgaaa | tgttttgata | ctttttgata | caatttgcta | ataactgttt | tgtagaatgc | 360 |
| ctgccgggggt | tttccacctc | atccctttcc | tcc | | | 393 |
| <210> 574 | <211> 397 | <212> DNA | <213> Homo sapien | | | |
| gcacgaggct | gccccggagct | gcctggggttg | cgctgccggc | cacgtccccg | cgccgggcct | 60 |
| caggctcctt | cctactgtcc | gagggccacc | aggccgccgg | gggcctgctg | cgcccggtatg | 120 |
| cgctctgttac | tagagtggag | agtctacctt | cgtctcacat | gtgccacaaa | ggatggcatg | 180 |
| gccccgggagt | gccccaccac | gtggctttca | ccccctgcaa | agccagactt | cgcccgagcga | 240 |
| cacagtgtca | agccacagc | tctccaagga | ggaagatggg | ccaggctggg | agcatccct | 300 |
| tagcagcagc | ctctgatccc | ttggccaagc | aggagggaaac | cattagcagc | ctgaggagct | 360 |
| ggctggctgg | gagcctcggn | gaccgcccag | ccttgct | | | 397 |
| <210> 575 | <211> 397 | <212> DNA | <213> Homo sapien | | | |
| ccatcgatt | cgaattcggc | acgaggctta | gggaacagga | gtgaacagac | ttcagcccca | 60 |
| cctggcaggg | gctggctccc | gaggttgggc | ccagtccctg | agggtctgct | ctgctacggg | 120 |
| tctgccccttg | agtggccttc | cgtggagggt | gtgtgaccag | gtggatgggtg | cagggcctct | 180 |
| ggagcccctct | cctcaggagc | agtcctcagc | ctttttctgt | aaaagacttt | tcttttggtgt | 240 |
| tctaggtggg | cagcagggttc | caggctgggtg | tttacaatct | cggaggaagt | gcgatgggtt | 300 |
| ctgttctttt | gacagttcag | tctgatttca | agtcagtcga | aagcgaacca | gaagcacggg | 360 |
| gcacagcagc | tcctctggct | gtgtagacag | acctggn | | | 397 |
| <210> 576 | <211> 394 | <212> DNA | <213> Homo sapien | | | |
| ggcacgaggg | tagggctgtg | ctgcgcggctc | cttcccatte | accctagtct | ggcgctcgcc | 60 |
| ggcgtggggc | ggccggacct | tcgccgcttc | caggaagggc | cacaacggcc | gtcggaccac | 120 |
| ggcgcgggcg | ccagttcctt | tatagttttg | ttcagaaaaa | catatggaga | cgtttatacc | 180 |
| cattgatttg | acaactgaaa | atcaagagat | ggacaaggag | gaaaccaaga | caaaaccaag | 240 |
| acttttaaga | tatgaagaga | aaaaatatga | agatgtgaaa | ccattagagt | ctcaaccagc | 300 |
| tgaatatagca | gaaaaggaaa | cattggaata | taaaacaagt | agaacaatct | ctggatcttt | 360 |
| tgaagcngag | gaaaccggag | gattacctta | gaga | | | 394 |
| <210> 577 | <211> 386 | <212> DNA | <213> Homo sapien | | | |
| ggcacgaggg | gaagtgccag | gaagaggagg | gtggccatgc | ctggccattt | cctgatacct | 60 |
| gtgctagtga | cggccgcggg | gtgtccactg | gaaagaaaca | ctggcgtgca | cggtgtgac | 120 |
| tgtggtttca | gcagttctga | gacaagagcc | ttccaagtgc | ggggctgggg | agcagagtgc | 180 |
| gggagctcct | gagtcctggg | ggcctccgcg | cctcacagca | tgggcacatg | tgggacagaa | 240 |
| ggcctaattg | ggtgcctgag | ggtggcctgg | ttgctgtccc | cccagggtgg | gaccatgagc | 300 |
| gagtgggggt | ggcacacggg | ctcagctctc | tgtggccggg | gtggctcctc | ttgcggggact | 360 |
| caacgtcagc | cccaaggcga | tgttca | | | | 386 |
| <210> 578 | <211> 386 | <212> DNA | <213> Homo sapien | | | |
| ggcacgaggg | ctcctggaaa | tgaagatgag | ctccacctgg | cacccgagcn | nnttggtgtg | 60 |
| ccccctcac | tgagggggcc | ccccgcaccc | gggaggagac | gcgggacttg | gtccacgctc | 120 |
| cgttaccctt | gacctggaaa | cgctcgagcc | tgtgtggtga | ggagcagggg | tcccccgagg | 180 |
| aactgaggca | gcgggagggc | gctgagcccc | tgggtggggc | ggtgcttctt | gtgggtgagg | 240 |
| caggcctgcc | ctggaaacttt | gggcctttgt | ccaagccccg | gcgggaactg | cgacgagcca | 300 |
| gccccgggat | gattgatgtc | cggaaaaacc | ccctgtaagc | cctcggggca | gaccctgcct | 360 |
| tggaggggaga | ctccgagcct | gctgaa | | | | 386 |
| <210> 579 | <211> 386 | <212> DNA | <213> Homo sapien | | | |
| ggcacgagga | gagagagaga | gagagagaga | gagagagaga | gagagagaga | gagagagaga | 60 |
| gagagagaga | gagagagaga | gagagagaga | gagagagaga | gagagagagt | cttttttttt | 120 |
| ttctctacct | ataaaaaacc | cccccgctgc | gtgtgtgtgg | ggggggacac | ccagaaaaca | 180 |
| cactatattc | tctctctctc | tgggcgcgcg | agagagagca | cacacggggg | ggaggggaga | 240 |
| aagcacgctc | tccccccccc | ccgtgttttt | tttttttttt | ttggcccccc | cccaacaaaa | 300 |
| aaaccacctt | tggtttcccc | ccccctccgg | gagaacaagc | cctttccccc | tttccatta | 360 |
| aaacagccct | tccccccccc | ccccct | | | | 386 |
| <210> 580 | <211> 399 | <212> DNA | <213> Homo sapien | | | |
| gattcgaatt | cggcacgagc | tcacaccaca | gctgagaggg | aaaggaaggt | tggatggcg | 60 |
| gatcgccaag | cgcgccccca | cctctcctgt | ggtactgggg | tccctaaagc | cgacccccgc | 120 |
| tccggcgggg | ctcgccggcc | cccaagtcgc | cagccgctta | cctcacaatc | ccgcttggac | 180 |

CLAIMS

We claim:

1. A library of polynucleotides, the library comprising the sequence information of at least one of SEQ ID NO:1-3351.
2. The library of claim 1, wherein the library is provided on a nucleic acid array.
3. The library of claim 1, wherein the library is provided in a computer-readable format.
4. The library of claim 1, wherein the library comprises a polynucleotide corresponding to a gene differentially expressed in a cancer cell of high metastatic potential relative to a control cell, wherein the control cell is a normal cell or a cell of low metastatic potential, wherein the expression is greater in the metastatic tissue, and wherein the sequence is selected from the group consisting of SEQ ID NOs:14, 137, 151, 152, 171, 200, 254, 262, 271, 348, 412, 472, 507, 520, 530, 588, 623, 637, 660, 678, 680, 700, 714, 774, 812, 834, 901, 937, 976, 1168, 1333, 1352, 1520, 1524, 1546, 1550, 1574, 1580, 1590, 1599, 1607, 1622, 1706, 1752, 1768, 1769, 1780, 1781, 1799, 1803, 1811, 1851, 1856, 1867, 1872, 1875, 1884, 1919, 1923, 1939, 1975, 2024, 2045, 2060, 2071, 2118, 2119, 2128, 2135, 2177, 2181, 2184, 2185, 2190, 2193, 2232, 2239, 2283, 2311, 2314, 2338, 2378, 2393, 2394, 2395, 2398, 2460, 2490, 2505, 2514, 2540, 2542, 2597, 2607, 2640, 2657, 2669, 2670, 2674, 2679, 2684, 2707, 2724, 2757, 2776, 2804, 2818, 2906, 2959, 2964, 2968, 2976, 2980, 2987, 3010, 3043, 3047, 3050, 3071, 3072, 3092, 3095, 3097, 3140, 3157, 3173, 3187, 3203, 3210, 3212, 3220, 3236, 3249, 3264, 3284, 3288, 3305, 3309, 3318, 3330, 3331, and 3335.
5. The library of claim 1, wherein the library comprises a polynucleotide corresponding to a gene differentially expressed in normal colon tissue relative to colon cancer tissue, wherein the expression is greater in the cancer tissue, and wherein the sequence is selected from the group consisting of SEQ ID NOs:7, 164, 734, 836, 928, 965, 987, 1026, 1044, 1119, 1226, 1227, 1251, 1316, 1429, 1442, 1540, 1553, 1560, 1577, 1588, 1610, 1620, 1626, 1673, 2416, 2749, 2976, 3129 and 3132.

541

6. The library of claim 1, wherein the library comprises a polynucleotide corresponding to a gene differentially expressed in normal colon tissue relative to colon cancer tissue, wherein the expression is greater in normal tissue than cancer tissue, and wherein the sequence is selected from the group consisting of SEQ ID NOs:105, 198, 465, 489, 745, 859, 976, 1011, 1045, 1138, 1226, 1251, 1253, 1392, 1474, 1559, 1571, 1589, 1591, 1607, 1608, 1643, 1753, 1764, 1766, 1782, 1811, 2749, 2784, 2790, 2805, 2976, 3128, 3129, 3146, 3150, and 3151.

7. The library of claim 1, wherein the library comprises a polynucleotide corresponding to a gene differentially expressed in normal human prostate cells relative to human prostate cancer cells, wherein the expression is greater in normal cells than cancer cells, and wherein the sequence is selected from the group consisting of SEQ ID NOs:53, 446, 1410, 1754, 1801, 1845, 2060, 2143, 2632, 2899, and 3338.

8. The library of claim 1, wherein the library comprises a polynucleotide corresponding to a gene differentially expressed in normal human prostate cells relative to human prostate cancer cells, wherein the expression is greater in cancer cells than normal cells, and wherein the sequence is selected from the group consisting of SEQ ID NOs:86, 93, 687, 1269, 1581, 1647, 1649, 1710, 1717, 1772, 1960, 2987, 3128, 3132, 3150, 3222, and 3268.

9. An isolated polynucleotide comprising a nucleotide sequence having at least 90% sequence identity to an identifying sequence of SEQ ID NOs:1-3351 or a degenerate variant or fragment thereof.

10. A recombinant host cell containing the polynucleotide of claim 9.
11. An isolated polypeptide encoded by the polynucleotide of claim 9.
12. An antibody that specifically binds a polypeptide of claim 11.
13. A vector comprising the polynucleotide of claim 9.
14. A method of detecting differentially expressed genes correlated with a cancerous state of a mammalian cell, the method comprising the step of:

detecting at least one differentially expressed gene product in a test sample derived from a cell suspected of being cancerous, wherein the gene product is encoded by a